

# ALGEBRA

#### LECTURES DELIVERED TO POST-GRADUATE STUDENTS OF CALCUTTA UNIVERSITY

BY

# FRIEDRICH WILHELM LEVI, DR.PHIL.NAT. HARDINGE PROFESSOR

PART III—CONTINUED FRACTIONS
PART IV—APPROXIMATE SOLUTION
PART V—MATRICES RESULTANTS

Parts 311 - 51



5114

PUBLISHED BY THE
UNIVERSITY OF CALCUITA
1937

0

BCU 1738

#### PRINTED IN ENDIN

PRINTED AND PUBLISHED BY EMPEROPRALAL CANCELLY.

109 947

Seg. No. 1062B-December, 1997-c

#### PREFACE "

The third and the fourth part of these lectures deal with some classical portions of Algebra. Obviously a strict selection has been necessary, and the author gave preference to those portions of Continued Fractions (Part 111) which are connected with the theory of numbers. In Part IV (Approximation solution) much importance has been given to how to carry out the calculations. Of course this branch of Algebra has to be considered from quite a different point of view from that of General Algebra. By approximation numerical values should be found out in the quickest and castest way; so the reader cannot get a real insight into the nature and importance of the different methods without knowing the difficulties a practical reckener has to face. The hints given by the author for the simplification of calculations do not involve the use of slide rules, calculating machines and graphical methods as these expedients are not familiar to our students. In this as well in some other items a later edition of these lectures is expected to show some alteration.

Part V contains the most important theorems on matrices with application on Hermitian and quadratic forms. Here the student may have the antisfaction of seeing in a nutshell the greater part of what he learned on Analytic Geometry.

As in the preferes of Part I and of Part II, I have much pleasure in delivering my thanks here to friends and kind collaborators. The Vice-Chancellor of our University, Syamaprasad Mookerjee, Esq., M.A., B.L., M.L.A., Barrister at Law, gave me every possible help to get these lectures printed in a very short time, and is fully entitled to the grateful thanks of the author as well of the students. Proofs and manuscripts have been carefully revised by Mr. A. C. Chowdhury, M.Se., Research Scholar of Calcutta University. The Calcutta University Press kindly co-operated in carrying out the printing in the proposed time.

F. W. LEYI

CALGUTTA, ASUTOSH BUILDING, October, 1937



## PART III. CONTINUED FRACTIONS

£1.	CANGUAL PROPERTIES OF CONTINUED PRACTIONS	1,
	<ol> <li>Principal definitions. 2. The h. c. f. 3. Proper and improper equivalence. 4. Periodicity.</li> </ol>	
f:2.	REPRESENTATION OF THE POSITIVE NUMBERS BY CONTINUED	В
	1. Unique representation of an irrational number by a continued fraction. 2. Interpretation of the continued fractions as positive numbers. 3. Distribution on the real axis, approximation by convergents.	
9 15.	PRESONG CONTINUED FRACTIONS WITH INTEGRAL CORFFICIENTS	14
	1. Representation of quadratic numbers. 2. Purely periodic continued fractions. 3. Approximation of quadratic numbers. 4. Reduced quadratic numbers. 3. Square roots.	
6.4.	APPLICATION OF THEORY OF NUMBERS	21
	1. Linear equations. 2. Pell's aquation.	
§ 5.	Continued practions whose elements are $\phi(x)$	22
	<ol> <li>The field of the elements φ(x).</li> <li>Representation of φ(x)</li> <li>by a continued fraction.</li> <li>Approximation of φ(x)</li> <li>Interpretation of continued fractions as elements φ(x).</li> </ol>	
₫ 6.	CONTINUED FRACTIONS WITH RATIONAL RESERVES	29
	1. Convergence. 2. Test of divergence. 3. Test of convergence. 4. Test of irrationality.	

#### INDEX

	PART IV. APPROXIMATE SOLUTION	Page
B 1.	Housen's scheme "	3.5
12.	I'm noors of a real and complex roots. 2. Budan-Fourier's theorem. 3. Sterm's theorem. 4. Legendre's polynomials. 3. Systematical investigation of the real roots. 6. Regula faisi and Newton's method. 7. Poulain's theorem.	42
§ 3.	GRAPPE'S METHOD	55
g 4.	ROOTS OF COMPLEX POLYNOMIALS	0.9
g ö.	Interpolation	64
	PART V. MATRICES. RESULTANTS	
9 1.	L Fundamental definitions. 2. Ring of Matrices.  8. Vectors.	71
( 2.	TRANSPORMATION OF A 18TO A ROBBAL-FORM  1. Vectors with invariant direction. 2. Polynomials of a matrix; the characteristic polynomial. 3. Case when the roots are different. 4. Case of one multiple root. 5. The invariant vectorspaces corresponding to the different roots.  6. The normal-form.	76
§ 3.	Sour enogenties of the normal-form and of the characteristic formountail.  1. Admissible bases for a normal-form. 2. Polynomials of which A is a root. 3. Linear substitutions of a complex variable.	88

vii

### INDEX

		Page
64.	THEORY OF MLEMENTARY DIVISORS	02
	<ol> <li>Matrices over S. 2. Congruent matrices. 3. Elementary divisors. 4. Normal-form of a matrix over S.</li> <li>Condition for congruence. 6. Connection between the normal-form of the matrix A over the field K, and the normal-form of A-zE over the ring K[z].</li> </ol>	
66.	HERMITIAN AND UNITARY MATRICES. HERMITIAN AND QUADRATIC	
	FORMS will the	100
	1. Notions and notations. 2. A restriction. 3. Transfer- mation of a Hermitian matrix into its normal-form. 4. Hermitian and quadratic forms. 5. Fundamental theorem of real quadratic forms. 5. Geometrical inter- pretation.	
5-6.	RESULTANTS	10
	1. Resultant as a function of the roots. 2. Resultant as a determinant. 3. General theorem. 4. Linear dependence of the resultant on its polynomials. 5. Elimination.	
Com	RECTIONS TO PARTE I-V	11



# PART III CONTINUED FRACTIONS

#### 1 L. GENERAL PROPERTIES OF CONTINUED PRACTICES.

Let K be a field, S a subring of K, and A be a subset of K with the [1/1] following property: If a class of residues ‡ (0) of S contains an element of A, this class contains also the inverse of an element of A. Hence if

$$a_1a_1^2, a_1^{22}, \dots, a_{2n-1}, a_{2n-1}$$
 (2, 1)

denote elements of A and

$$a_1, a_2', a_2', \dots, a_1, a_2, \dots$$
 (1, 2)

denote elements of S, then every element of A can be represented either by

$$a = t + 1 + a^t \tag{1.8}$$

(1, 39)

or by a=r.

If s, g. K is the field of the real numbers. S the ring of the integers, and A the set of the real numbers >1, the representation (1,3), (1,3') is always possible and s is uniquely defined by a as the greatest integer  $\leq s$ .

But the conditions (1,3), (1,3) hold also for other sats A in fields K, so an investigation in general terms is helpful.

then 
$$a_1 = x_1 + \frac{1}{x_2 + \frac{1}{x_3 + \dots}}$$
 (1, 5)

The representation of a by (5) is said to be a continued fraction. The formulae (1.4) can be continued till an element a, will be an element of S. If there is an m so that a = s\_ then the continued fraction is finite, otherwise it is infinite. If a can be represented by a finite continued fraction.

N.

it belongs to the quotient field Q of S, and every finite set of elements

$$x_{1}, \dots, x_{n}$$
 (1, 6)

of S defines an element of Q by the belp of (1,5).

We now define other sequences of elements of S by the following formulas.

$$P_{-1} = 0$$
,  $P_{+} = 1$ ,  $P_{+} = s_{+}P_{+-1} + P_{+-2}$ ,  $h = 1, 2, ...$  (1, 7)

$$Q_{-1}=1$$
,  $Q_{+}=0$ ,  $Q_{+}=s_{1}Q_{s-1}+Q_{s-2}$ , (1, 8)

$$D_{k} = \begin{vmatrix} P_{k} & P_{k-1} \\ Q_{k} & Q_{k-1} \end{vmatrix}$$
 (I. 0)

then from (1.7) (1.8) (1.9) it follows that

$$D_1 = -D_{k-1}$$
, and so  $D_s = 1$ ,  $D_1 = (-1)^k$ . (1, 10)

From (1.9) and (1.10) it follows that P2 and Q2 have no other common factors in S than unities and that

$$P_{x}: Q_{x} = P_{x-1}: Q_{x-1} \approx (-1)^{x}: (Q_{x}Q_{x-1}).$$
 (1, 11)

Lat a be an arbitrary element # 0 of K, then we get a uniquely defined sequence of elements a, a, ..., a, g by the equations

$$a_i \approx a_i : a_{i+1}$$
 (1, 12)

From (1,4) we get by multiplying the equations with a<sub>4</sub>, o<sub>5</sub>, ... respectively

$$a_i = a_1 a_{i+1} + a_{i+0}, \qquad i = 1, ..., n$$
 (1.49)

From (1.4') (1.7) (1.8) we get

$$P_3a_{k+1} + P_{k-1}a_{k+2} = P_{k-1}(a_1a_{k+1} + a_{k+2}) + P_{k-2}a_{k+1}$$
  
=  $P_{k-1}a_k + P_{k-2}a_{k+2}$ .

By the repeated application of this formula we see that for ick

$$\mathbf{P}_{1}a_{1+1} + \mathbf{P}_{1-1}a_{2+2} = \mathbf{P}_{1}a_{1+1} + \mathbf{P}_{1-1}a_{1+2} = \mathbf{P}_{1}a_{1} + \mathbf{P}_{-1}a_{2} = a_{1}, \ (1, 13)$$

and by making a similar calculation with the elements Q we get

$$Q_{i}a_{i+1} + Q_{i-1}a_{i+2} = Q_{i}a_{i+1} + Q_{i-1}a_{i+2} = Q_{i}a_{1} + Q_{-1}a_{2} = a_{2}. (1, 18)$$

Hence

$$(P_{1} x_{k+1} + P_{k-1}): (Q_{1} x_{k+1} + Q_{k-1})$$
  
=  $(P_{2} x_{k+1} + P_{k-1}): (Q_{2} x_{k+1} + Q_{k-1}) = x_{3}$ . (1, 18°)

If we multiply (1.13) with  $Q_{x-1}$  and (1.13) with  $-P_{x-1}$  and add, then it follows from (1.10) that

$$(-1)^{\delta} a_{k+1} = a_1 Q_{k+1} - a_2 P_{k+1},$$
 (1, 14)

From (1,12), (1,13") and (1,10) it follows that

$$a = \frac{P_d}{Q_s} = \frac{(-1)^{1-1}}{u_H} = \frac{a_{1+1}}{Q_1}. \tag{1.15}$$

Hence if the continued fraction is finite, as any = 0

$$a_1 = \frac{P_a}{Q_a}$$
, (1, 15)

If (1, 4) holds, the elements  $s_1, s_2, ..., s_n$  are said to be the elements of the continued fraction; the quotients  $P_1:Q_n$  are the convergents and  $a_{n+1}$  is called a complete fraction. As  $a_1$  is uniquely defined by these elements, we shall denote  $a_1$ , if  $a_{n+1}$  exists, by

$$a_1 = (a_1, \dots, a_n \mid a_{n+1}) = \frac{P_n a_{n+1} + P_{n-1}}{Q_n a_{n+1} + Q_{n-1}}$$
, (1, 5')

and if s, is the last element of the continued fraction.

$$a_3 = (s_1, ..., s_s) = \frac{P_s}{Q_s};$$
(1, 3")

from (1.4), (1.57), (1.57) it follows that for  $k \le n$ 

$$a_{\alpha} = (s_{\alpha}, \dots, s_{\alpha}) \circ_{\alpha \in \Sigma}$$
, respectively (1. 5.7)

$$w_{k} \approx (x_{k}, ..., x_{k}),$$
 (1, 5,7)

Let P', and Q', be defined by

$$P'_{-1} = 0$$
  $P'_{+} = 1$   $P'_{+} = s_{i+i-1}P_{i-1} + P_{i-2}$   
 $Q'_{-1} = 1$   $Q'_{+} = 0$   $Q'_{+} = s_{i+i-1}Q_{i-1} + Q_{i-2}$  (1.16)

then the following formulae hold:

$$\begin{vmatrix} P'_{i} & P'_{i-1} \\ Q'_{i} & Q'_{i-1} \end{vmatrix} = (-1)^{i}$$

$$a_{i} = P_{i}a_{i+1} + P'_{i-1}a_{i+i+1}$$

$$a_{i+1} = Q'_{i}a_{i+1} + Q'_{i-1}a_{i+i+1}$$

$$(-1)^{i}a_{i+1} = a_{i}Q'_{i+1} - a_{i+i+1}P'_{i+1}$$

$$(1, 17)$$

4

MADLERY

$$a_1 = \frac{P}{Q} \qquad \qquad P = \frac{P}{Q$$

Frems

$$P_{n-} = s_n P_{n+1} + P_{n-2}$$
  $Q_n^* = s_n Q_{n-1} + Q_{n-2}$   
 $P_{n+1} = s_{n-1} P_{n-2} + P_n$   $Q_{n-1} = s_{n-1} Q_{n-2} + Q_n$ 

$$P_1 = x_1 - x_2$$
  $Q_2 = x_2$   
 $P_n = 1$   $Q_1 = 1$ 

we get the representation (P. P., ) and of Q.  $Q_{n-1}$  as finite sent need fractions. Then

$$\frac{P_n}{P_{n+1}} = (s_n, \dots, s_n, s_1) - \frac{Q_n}{Q_{n+1}} = (s_n, \dots, s_n)$$
 (1, 17a)

[1 2] There is a very close connects in between the fin to continued fraction and and the algorithmus of the for f.

Let  $a_1$  be represented by a finite continued fraction  $a_1 = a_2 = a_3 = a_4 = a_4$ . Hence  $a_1 = a_2 = a_4 = a_4$ . Hence  $a_2 = a_4 = a_4 + 0$  therefore  $a_{1,2} = 0$ . From (1.17), (1.13), (1.14), (1.14), an get therefore

$$a_1 = P_n a_{n+1}, \quad a_2 = Q_n a_{n+1}$$

$$(-1)^n a_{n+1} = a_1 Q_{n+1} - a_2 P_{n+1}.$$
(1, 16)

Hence at a P. Q. be ange to the quotionthead of S

Fivery common factor for and  $r_0$  is a factor of  $r_{-1}$  and  $r_{-1}$  is a common fact radio, and  $r_0$ . Hence  $a_1$  and  $a_n$  have an h of and this can be represented interry by  $a_1$  and  $r_0$ . In a representation by two relative by prime commutator b, as

$$P_n Q_{n+1} - Q_n P_{n+1} = (-1)^n$$
,

Let every element of the quotients id at 8 be represented in try a finite continued fraction and let \* \*" \$0 be two arbitrary elements of 8 then \* \* => can be represented by two Ye streety principles of the at the financial beautiful to a linear and homogeneous manner by those elements.

Hence a' = p : q and pq' + qp' = 1.



Therefore an ar and appear if a area = a

Hence s'=ps = sp = p = s So the art transecements s = fS have an hor f (s' sp = sm libs is represented in a pear out hour group transection by s' and s' from this count force in it f loss that it is not pear the to represent the motions of the elements it an archivers of a by timber on the motion as not un one there in sat to a quit out if two comments of S who he can be represented by an initiate continued fraction the finite continued fraction the finite continued fractions do not from a total in these cases.

Late stunction N(s) which takes practive integral value only to defined for every element of the first new territorial statements at the thorough two other elements so and a section.

$$e^{\pm a_1 e^r + e^{rr}}$$
 and that 
$$e^{rr} = 0 \quad \text{or} \quad \mathcal{N}(e^{rr}) < \mathcal{N}(e^r), \qquad (1, 19)$$

Then 
$$+ x = x_1 + x' + y = x_1 + x + \cdots$$
 and as  $N(x') > N(x'') > N(x'') > 0$ 

nither

are all integra numbers, the sequence of a most be finds force as at is a finite continue I function. From these considerations we get the following theorem.

The rest Let S be an ext f containing 1 and et ap abve integer S a) be defined satisfying the conditions (1.19 for every channel of S than the quot entheld of S we be form d by the finite contained fractions of S, and the highest common factor of two elements of, and of S or given by and in (1.18) where P<sub>a=1</sub> Q, have the against a given by (1.6), (1.7), (1.6).

1 at A and B be the matro es of solute aliens of the intel

$$\begin{split} \mathbf{A} &= \begin{pmatrix} a_1 & a_2 \\ a_3 & a_4 \end{pmatrix} & \mathbf{B} & \begin{pmatrix} b_1 & b_2 \\ a_3 & b_4 \end{pmatrix} & \det \mathbf{A} = a = \pm 1 \\ & \det \mathbf{B} = a = \pm 1 \\ & \det \mathbf{B} = a = \pm 1 \\ & \det \mathbf{B} = a = \pm 1 \\ & \det \mathbf{B} = a = \pm 1 \end{split}$$

$$\mathbf{A} &= \begin{pmatrix} a_1 & a_2 \\ a_3 & a_4 \end{pmatrix} & \mathbf{B}_1 &= \begin{pmatrix} b_1 & b_2 \\ a_3 & b_4 \end{pmatrix} & \det \mathbf{B} = a = \pm 1 \\ & \det \mathbf{B}$$

che be transfermed by hinter, and to be transformed by Bunton, then it follows

The product of two matrices has been debut in first 1, ; 27 m 1 st has then shown in . 38 that the determinant of a product of matrice at a copy of to the product of the determinants. The second with his we have to exhaust other an early a president at other n

An element of K a shall to be equilibrially and a get 1 by transforms the equal to a reflections substitute a with intermediate 1. He is all (a) and 3 the laws that these quarterese satisficable can be easily of reflex vity symmetry and transformly see that to {1 3 } and therefore that equal transform of K ato causes so that two consists of K are equivalent fund only (they teleng to the same easily).

By 
$$\binom{n+1}{1} = \binom{-1}{0}$$
 the charact there are transferred at  $r$  is ness all

elements I S are equivalent. From (1,1%, and 1% to Lowe that the ments is any left need by 1% are as equivalent 1% who purely course that e attended fraction by 1% are as equivalent to the example of the computation of the course of the computation.

the same there exists the best rise det A = 1 in the second case de A = 1 to det A = d A in 1,20 h lie the set in of proper equipment as well as the pet in decayed the same kind the set in of proper equipment as a symmetre into By combining two equipments of the same kind we get a proper upon the area on 3 by combining two equipments of the same kind we get an improper equipment a property contact the deal, for the tractic E has the letter mantal. If in a class of equipment a make improperly equipment to itself, if a a transformed into a 1 y 1, and det k = -1 then an arbitrary ordered defined these matrices has the determinant 1, the other has the determinant 1 so each element of the class in properly and improperly equipment to a and therefore every is ment a at the same time properly and improperly equipment to every chiefer every is next as the same time properly and improperly equipment.



then a become transformed into the BA, where dot A = 1, dot B = 1, then a become transformed into the BA, where set BA = -1, and therefore is emprepariting valent to the factor of the dose compressively and improper vacquivalent to every therefore in the factor of the dose of the factor of the f

#### Honco

The m In a class of equivalent elements, either every chiment is properly at 1 major perly equivalent to the views therefore each test related to the end of the end

As I a transfirmed out stood by \$ \$1 a very past of elements a property and improperty equivalent in the class containing the finite contained frontions

Let contransf receil into itself by A. Itan

$$a = \frac{c_1 a - a_2}{a_1 a_2 a_3}$$
 holds, better  $a_2 a_3 + (a_2 - a_3) a_1 a_2 = 0$ . (1, 21)

. There are 3 different cases

I 
$$a_1 = a_2 + a_3 = 0$$
 In this case  $A = L$  or  $-1$ 

By these to nefermations every element is transferous domain test. It and -E generale proper equivalences.

- 2 (1, 21) is a reductive patrons and a line is possible only if the an element of the quotientfield Q of B
- A (1, 21) is reclucible. In this case a single rate to Q and fielder 2 over Q.

I rest these considerations it for own that not every extremt in improperly equivalent to itself.

Let the elements a , a undefined by 1 t, be not ad different by but

$$a = a$$
 then

$$\alpha = \sigma_{-1}$$
 ,  $\sigma_{-1} = \sigma_{-1}$  ,  $\sigma_{-1} = \sigma_{-1}$  (a)



8

#### ALGEBRA

Hence a ran be represented by an antinite period c continued fraction with

Pransfermed at itself by the matrix 
$$\mathbf{D} = \begin{pmatrix} \mathbf{P}_{t}^{i} & \mathbf{P}_{t+1}^{i} \\ \mathbf{Q}_{t}^{i} & \mathbf{Q}_{t+1}^{i} \end{pmatrix}$$
 where

del 15 - 1 and becomes therefore equivalent to self by the transforms

tron From (1, 13°) it follows that 
$$a_1 = \frac{P_{-1} + \frac{P_{1-2}}{Q_{1-2}}}{Q_{1-2} + \frac{Q_{1-2}}{Q_{1-2}}}$$
 belongs to the field

Q a J. Hence at and at each equation of agree 2 with configuration S. The same holds for eq. eq. :

[1/4] Let now a periodic acquence a , v., a1 a.,

of elements of Silvegore. It is not errors that an any extension of the quotient field Q of Sithers exists an element which can be represented by the intinite periods and a sed fraction.

$$(a_1, ..., a_m, a_1, ..., a_m, ...)$$
,  $(1, 22)$ 

If such an element exists at is a C certain that it a clonent a uniquely before a the field. But if there is a field in which there exists one and only one element a represented by the period c continued fraction it all then

$$a = (a_1, ..., a_m, a_1, ..., a_m, ...) = (a_1, ..., a_m \mid n)$$

hat is and the is the case considered past before

- 2. REPRESENTATION OF THE POSITION NUMBERS BY CONTINUED PRACTICAL
- [2/1] Let the elements ( 1 -f the set A be the real numbers  $\geq$  1

  not et S be the ring of the integers | then the representation by the form due
  (1, 3) and (1, 3)

$$a = a+1$$
:  $a'$  or  $a = a$ 

is always possible. If a ir not an integral number there is

$$1 \le s \le s \le s + 1, \ s' = 1 \cdot ss + 1 - st$$

and his re-resert, on a night fair fair in a graft posts thinners,

 $a=a'+1,\ a'\approx 0,\ 1,\ 2,\ \dots$  , there exist two possibilities

Dispersion to the present in the first of the property of the present of the pres

In process of all a constants of the constants of the constant of the constant

$$n_1 \simeq (n'+1)$$

$$n_1 \simeq n'+1$$

nuction of the continued freehom t

and

If paret plager, a unit yiefn a process of an entered in minimpely deposed of Anthronou pely deposed of Anthronou pely deposed in the formal of the example of a process of the example of

$$u_{r+1} = (r'+1)$$
 , or  $u_{r+1} = r'$  ,  $u_{r+1} = 1$ ,

and there exc t two and only two representations of a

$$a_1 = (a_1, \dots, a_r, a_r', 1),$$
 and  $a_1 = (a_1, \dots, a_r, a_r', 1).$ 

ed to refer the second to the second the second terms of the

$$H \mapsto H \leqslant 1 \quad \stackrel{f}{\longrightarrow} \quad \Rightarrow 1, \text{ and } \quad \stackrel{f}{\longrightarrow} \quad 1 \quad 0, \text{ are called}$$

derations is given by the forming, the norm

fenction and affing the condition and of the minimum or to the other than



representation is unique and the fint of all feating sinflate. If the aims because on there exists a single continue of the aims and feating and the there's an absolute to the continue of the fint of the first of the

[22] We write the reserve the region to the every sequence set by no 1 1 form and on your contract.

Let spend on the an infinite sequence sate fying the conditions (I 10). The near case P and Q should be being as 1 7) and (I 8).

$$P_{-1} = 0$$
,  $P_0 = 1$ ,  $P_1 = s_1$ ,  $P_2 = s_2$ ,  $P_3 = s_3$ ,  $P_4 = s_4$ ,  $P_4 = s_4$ ,  $P_{-2}$ ,  $P_{-3} = 1$ ,  $Q_0 = 0$ ,  $Q_1 = 1$ ,  $Q_2 = s_3$ ;  $Q_4 = s_4$ ,  $Q_{-1} = 0$ .

From Park 21 to the gammathemate distance but

$$0 \le P_1 < P_0 < P_3, ... < ..$$

$$Q_0 = 0 < Q_1 \le Q_0 < Q_3, ... < ... \text{ bold}$$
(2, 8)

The quarter  $\frac{P_{p}}{Q_{p}}$  are for k < n the convergence of a character to

$$C_1 = A_0 = \frac{P_0}{Q_0} = SC_1 = 0$$
 therefore apply 1-15 no.  $\approx \frac{P_0}{Q_0}$ 

Hen e

$$\begin{array}{lll} P_{i} & P_{i} & = & (1)^{k+1} a_{k+n} & \geq 0 & \text{if } k \geq 1 \text{ odd} \\ Q_{i} & Q_{i} & a_{i} & Q_{i} & < 0 & \text{if } k \geq 1 \text{ oven} \end{array}$$

Rende

$$\frac{P_{1}}{Q_{1}} = \frac{P_{2}}{Q_{1}} < ... < \frac{P_{2m+1}}{Q_{2m+1}}$$

$$\frac{1}{Q_{2}} > \frac{P_{3}}{Q_{4}} > ... > \frac{P_{3m}}{Q_{4m}} \quad \text{for } m = 1, 2, 3, ... .$$
(2.4)

The quantity  $\frac{P_{\rm q}}{Q_{\rm m}}$  form therefore two sequences of the series is necessary, the other is decreasing and every number of the first equal case when every

number of the second one the attrible 
$$\begin{pmatrix} P_+ & I_+ \\ Q_{n-1} & Q_n \end{pmatrix}$$
 form therefore a



set of infersals cach of them a not ded now good er pre alogat-

As 
$$\frac{P_{i}}{Q_{i}} = \frac{P_{i+1}}{Q_{i-1}} = \frac{(-1)^{2}}{Q_{i} Q_{i-1}}$$
 . [see (1,31.]

The length of them of reconstruction and the configuration of the config

$$\frac{\Gamma_2}{Q_{2m}} = a \times \frac{1}{Q_{2m}} + \frac{1}{2} = \frac{1}{2} + \frac{1}{2} + \frac{1}{2} = \frac{1}{2} = \frac{1}{2} + \frac{1}{2} = \frac{1}{2}$$

with problem of the state purpher which which the y-hi sequence  $x_1 x_2 = x_1 x_2 + x_2 x_3 + x_4 x_4 + x_5 x_5 +$ 

actuated within all the interval 
$$\begin{pmatrix} P_n & P_{n+1} \\ c_n & Q_{n+1} \end{pmatrix}$$
, but it may so happen that

there is no such number on heat the represent in it is a commend fraction furnishes about respect which is from the same recommend. Of course the complement near the same W. is allow the highest parameters. We saw about the highest parameter which highest parameters at the detribution of the communications on the name of the real numbers.

Figure 3 for P. Q. to the conversants of  $(s_1, s_2, ...)_i$  and  $\Gamma^{i,i} Q^i$  be the last convergent of  $s_1, ..., s_{n-1}, ..., s_{n-n-1}$ , where i = 1, ..., and i > 0.

Then

Proof. 
$$Q_{n+1} = s_{n+1} Q_n + Q_{n-1}$$
 
$$Q' = (s_n + t) Q_{n-1} + Q_{n-2} = Q_n + t Q_{n-1}$$
 
$$\frac{1}{t} Q = \frac{1}{t} Q_n + Q_{n-1} \leq Q_{n-1}$$

The equanty hade ply if n = 0 or if  $t = \epsilon_{n-1} = 1$ 

um-1

Applying (1, 11) we get

$$\frac{1}{Q_{n-1}} = \frac{1}{Q_n} = \frac{1 - Q_{n-2}}{1} = \frac{(-1)^{n+1}}{1} = \frac{1}{(-1)^{n+1}} = \frac{2.2}{1}$$

$$\frac{1}{Q} = \frac{P}{Q} \qquad \qquad \gamma = \left( \begin{array}{ccc} 1 & P & 1 \\ Q & Q & 1 \end{array} \right)$$

$$= \frac{(-1)^+}{Q'Q_{n+1}} = \frac{(-1)^{n+r_0}}{Q_n Q_{n+1}} = (-1)^{n+r_0} = \frac{Q' - Q_n}{Q'Q_n Q_{n+1}} = \frac{1}{Q_n Q} = -70$$

From 2 to an ear 10 years and 11 as the ty-

Arollon or the element of the end of the end

the operation of a sorte to emerce on afferent red national terresponds to the transfer of the contract of the transfer of the

1981 From the form owners had been used fractions are distributed on the rate of the rate name of the content to world write for on artificing finite net to a content.

$$(u_1, \dots, u_n) = [u_1 | t] \quad t = 0, 1, 2, \dots$$
 $u = [u_1, u_2] \quad u = 1, 2, \dots$ 

Let a be odd

$$[n \mid 0 < \lfloor n, n \mid < - < n, 2], < [n \mid 1] = [n \mid 1] < [n \mid 2] < .$$



If a wavet, the note is < in the region by > but the continued fractions having only an element, then we  $r_1 = r_1 > r_2 > r_3$  we get a fractions of the reason by the term red L > line attracts fractions  $\{r_1\}$  or begins  $n_s < 1$ ,  $r_1$  are a varied cute in  $r_2$  and  $r_1 + 1$ .

For the continued from a some way 4 to the same has been a some the present of the same and the same as the present of the same a

As every rational number is tope in class a total or distribution of the Connumber of metrics as the corresponding tope of the distribution of the classical contribution of the code panel contains the contribution of the code panel code in the contribution of the code panel code of the code

The super success the consideration of more really the following theorem.

Theorem. It s>0, and 
$$\frac{P_{q-1}}{Q_{2n-1}} < \frac{r}{s} < \frac{P_{q-1}}{Q_{2n}}$$
 then  $s > Q_{q-1} > Q_{q-1}$ .

Proof. From the supposition it follows directly that

$$0 = \frac{r}{r} = \frac{\mathcal{V}_{2r+1}}{\mathcal{Q}_{2r+1}} < \frac{\mathcal{V}_{2r}}{\mathbf{R}\mu_r} = \frac{\mathcal{V}_{2r+1}}{\mathcal{Q}_{2r+1}} = \frac{1}{4(2r^3(2r+1))}$$

and as a and Que-1 are positive

the middle part of the in you by is a integral patient wimber

ALGEBRA

the believe meanine but to be to promote a decret newser by the believe to a secure of the believe denominates are the present as he to express the transfer of the secure of the secure

#### 5.3. Pentodic continued priori a with internal co-residents.

3/1) I the beneath that for a non-deady non-section of the form and continued a selection of the form and continued a selection of the form and continued and and another than the form and and the form of the first transfer of the form of the form

where A By 2 top is referred by a periodic continued fraction

Proof. a less to be the root of polynomial

$$as^2 + 2bs + 0$$
; (8.1)

a the other hand is a resented by a cut naced fact is
a = (σ<sub>1</sub>, ..., σ<sub>n</sub>, λ). From (1, 13°) it follows that

If the  $x_1 P_n \leftarrow P_{n-1} = a_n x_1 P_n + P_n + (Q_n V + Q_{n-1} V + Q_n V + Q_{n-1})^2 = a_n V$  and  $x_1 = a_n V + (Q_n V + Q_n V + Q_n V + Q_n V)^2 = a_n V + (Q_n V + Q_n V + Q_n V + Q_n V)^2 = a_n V + (Q_n V + Q_n V + Q_n V + Q_n V + Q_n V)^2 = a_n V + (Q_n V + Q_n V)^2 = a_n V + (Q_n V + Q_n V)^2 = a_n V + (Q_n V + Q_n V + Q_n$ 

Let I be the second roll of a leand a to defer by

$$\beta = \frac{P_{n}\mu + P_{n-1}}{Q_{n}\mu + Q_{n-1}},$$
 hence 5.4)

 $\Delta \kappa = \frac{Q+\gamma}{Q+\gamma} = 0, \text{ to the effect of the proof of$ 



As 
$$\frac{P_{\alpha}}{Q_{\alpha}} = \frac{P_{\alpha-1}}{Q_{\alpha-1}} = \frac{P_{\alpha-1}}{Q_{\alpha}} = \frac{1}{Q_{\alpha}Q_{\alpha-1}} = \frac{1}{Q_$$

From (0, 2) it follows therefore

$$P = \begin{cases} Q_{n} & Q_{n} & Q_{n} \\ Q_{n} & Q_{n} & Q_{n} \\ Q_{n} & Q_{n} & Q_{n-1} & |B-n| > 2 \end{cases}$$

$$< 0 \text{ for } Q_{n} Q_{n-1} |B-n| > 2$$

Hen a little a certain integers the religible main regative therefore h = f + h independent in gratice. As  $h^2 = e^{-h} + h$  independent on a feete mancher of solutions h = h. If  $f = e^{h} + h = h$  in the constant be different annual fractions and solve  $f = e^{-h} + h$  in the first the same countries for the same countries for the same countries for the same countries.

Leb my 13, to be crist from my force to be a mellional for a series properly who became to make the force of a form 4 we know that the rest operate to a sect from a force or make the force of a form 4 and a so partie to be a force of a force on of a force or for every pure a force or interest or force or purely periods a number > 1.

A part to de quarter a equal mass of the rest of the 1 and the conjugate root sate  $x = 0 > \mu > -1$ . We will prove a with descriptoric y per odic continued fraction represents a reduct quadratic number.



16

#### A FORTURA

Let

$$a = (a_1, \dots, a_n) = (a_1, \dots, a_n \mid a)$$
 (8.7)

$$\xi = (a_1, \dots, a_k) = (a_1, \dots, a_k \mid \xi)$$
 (8.79)

be two purely contract if features, the elements at being the same in both continued feetures, but ordered in an inverse manner.

Lat. P., Q. be the convergents of a. Then (see 1, 175)

from c

Let 
$$\beta = \frac{-1}{2}$$
, then  $0 > \delta > -1$ , and

$$f_{-}(x) = Q_{n-2}x^{n} + (Q_{n-1} - P_{n})x - P_{n-1}$$

In , 
$$\frac{U - P}{U - C}$$
, also shoot  $f(x)$ , and as  $a > 3$ ,

the t is a different (b) one consider this atheretic

at one is the Lagrangian term of the Lagrangian to a, and  $\xi = -\beta^{-1}$ , then  $\xi$  is represented by (3, 7)

(very cottons from a quantity to the complete few times out to a possible or prostly credity and fraction; bence

Constary. Every quadratic number is equivalent to a reduced quadratic number

seminal files of a become r

$$a = \frac{a + \sqrt{D}}{b} = s + \frac{1}{2}$$
, where  $s < s < s + 1$ .

and a, b, s, D > 0 are integral numbers.



From these forms as we can find ut by a mosphere more ease he is the numbers of, it defaults obsputly the complete fractions a unit the numbers of a lamb to be a mosel fraction. As the not old fraction is period of one pure it in the period of and it is most on how to be stopped.

becomples.

1. 
$$a = \frac{1 + \sqrt{5}}{2}$$
 (barmonic section)  $D=5$ 

1.  $a = \frac{h}{2}$  (c)  $a = \frac{1 + \sqrt{3}}{2}$  (c)  $a = \frac{1 + \sqrt{3}}{2}$  (d)  $a = \frac{1 + \sqrt{3}}{2}$  (e)  $a = \frac{1 + \sqrt{3}}{2}$  (f)  $a = \frac{1 + \sqrt$ 

$$2c = \sqrt{26}$$
,

ony other dead,

hence a more I has example as very a nament for york and exact patentialion.

$\mathbf{P}_{\sigma} =$	1	Q,=	o
$\mathbb{P}_1 =$	5	$Q_1 =$	1
$P_{u} =$	,	V <sub>2</sub> ≈ 10	o
$P_{\rm total}$	913	$Q_{\pi} = -10$	L
$P_{ 0} =$	528 T	$Q_{\chi} = -105$	U
$\mathbf{F} =$	6-1-	Q <sub>5</sub> == [080]	ı
$\Gamma_{\rm ic} =$	sart 1	Q <sub>a</sub> = 104080	)

#### ALGEBRA

Hence  $a = \frac{5 \text{ and }}{100 \text{ gr}} = 0$  where  $0 < \epsilon < 10^{-13}$ . Therefore a = 5 totally 1s is the last two figures 1 mg un estad.

Voltholerror 
$$c = -c + \frac{P_+}{Q_+} < \frac{1}{Q_+Q_-} < \frac{1}{Q_+Q_-} < \frac{1}{c_+ - c_+Q_-^2}$$
 , it is useful to

at p tho ca culate in just before a sign and

Hence 2 2 . As P., Q. are increasing very abouty we will use another method

$$\sqrt{2} = \sqrt{2\pi}e^{-\frac{1}{2}} + 1$$
 If  $\sqrt{2} = \frac{P_{\infty}}{4e^{2}} + 1$   $\sqrt{2} = \frac{P_{\infty}}{10Q_{\infty}} \pm \frac{e^{-\frac{1}{2}}}{10}$ 

be we represent  $\sqrt{200}$  by a continue life too. D = 200  $-117 \pm 1$ 

Hence √200 =(14, 7, 28),

$$P_a = 1$$
  $Q = 0$ 
 $P_1 = 14$   $Q_1 = 1$ 
 $P_2 = 90$   $Q_2 = 7$ 
 $P_3 = 2786$   $Q_3 = 197$ 
 $P_4 = 19601$   $Q_4 = 1997$ 

$$\sqrt{200} = \frac{10001}{1000} = e$$
,  $0 < e < \frac{1}{Q_+ Q_+} < \frac{1}{28} \frac{1}{Q_+} < \frac{1}{2} \frac{10^{-8}}{Q_+}$ 

$$\sqrt{2} = \frac{1980 \text{ f}}{1.980} = e^{-1} = 0 < e < e < 10^{-15} = 1.514 \pm 1326 \text{ f}$$

true to eight figures after the decimal point

[3/6]



Exercises Prove that a  $2n = \sqrt{a^2 + 1}$ , and calculate  $\sqrt{2 + 1}$ ,  $\sqrt{82}$ .  $\sqrt{1 + 2}$ ,  $\sqrt{17}$ . Calculate  $\sqrt{3}$  due thy and also by hop of  $\sqrt{2 + 1}$ .

In order to prove the converse of the last the com, it is useful to consider the following lemma.

from a If a>1 and b<0 we complete quadrate numbers  $a=(a_1,a_1,a_2)$ , then all the complete fractions  $a_1=(a_1,a_2)$ ,  $a_2=(a_2,\ldots,a_{2n})$  are reduced a null are

 $P(\alpha, \beta) = \alpha = \alpha + \frac{1}{\alpha_1}, \ \beta = \alpha + \frac{1}{\beta I_1}, \ \alpha_1$  and  $\beta_1$  are conjugate numbers

 $a_1>0, \ \frac{1}{a_1}=a-\beta>a_{-\beta}$  bence the minduced, and its report in if

this procedure t follows that age are reduced

The rem Every reduced qualent continuous is represented by a purely periodic continued fraction.

fraction a s, s<sub>1</sub> s<sub>2</sub>) Let this number be reduced and of the period city of the continued fraction began with a only so let a first then it follows from the last comma that s s<sub>1</sub> , s<sub>2</sub> , is a reduced number to We will prove that the a impassible. Using the same notations as in the lemma we state

hence 
$$\beta_1 \Rightarrow \beta_{n+1}$$
  
 $\alpha \Rightarrow \alpha + \frac{1}{\alpha_1}, \quad \alpha_n = a_n + \frac{1}{\alpha_{n+1}}$   
hence  $\beta = a_n + \frac{1}{\alpha_1}$   
 $\beta_n = a_n + \frac{1}{\beta_{n+1}}$   
 $\beta_n = a_n + \frac{1}{\beta_{n+1}}$   
 $\beta_n = a_n + \beta_{n+1}$ 

but do a and  $a_{n-1}$  are reduced  $0 < -\beta < 1$  and  $0 < -\beta_{n-1} < 1$ 

hold, hence  $a = 1 < \frac{1}{\beta_1} < s$ , and  $s_n - 1 < \frac{-1}{\beta_{n-1}} < s$ . From  $\beta_1 = \beta_{n-1}$ 

st fellows therefore that a - r.

Theorem Lat. = V > 1 be restroyal then

$$n = (a, a_{\pm}, a_{\pm}, a_{\pm})$$
 (5), (8)

and 
$$a_1 = 2 < a_2 < a_3 = 2 < a_4 < a_4$$

hald literatures a sound 3,5 hald then at an armitona aguard root > 1

Prof. As > 1 months number < 0 months in a supposed to a set the manufacture of the frame of the

 $a_{\rm hol} = a = a = r + \frac{1}{r_1} + a_{\rm hol} = a_{\rm hol} = a_{\rm hol} = a_{\rm hol}$ 

$$a_1 = (a_1, \dots, a_n)$$
, (8.0)

Therefore the dates from the the remodel 2, that

$$-1:\beta_1=(s_0,...,s_1)$$
 (6, 9)

 $0 = \star_1 + t_1 + t_2 + t_1 + 1 + \dots + H(a|a-1| H_1 + \cdots + t_1)$ 

Therefore 
$$\{x_1, x_2\} = (2s, \frac{1}{s}, \frac{1}{s})$$
 (2s,  $\frac{1}{s}$ )

h le ( 1 H' is equivilent ( 1 10) conver y if a defined by (8, 8), (8, 8'), then (8, 10) holds

Let a pand a be defined by a uple (a), and let  $B = x + 1 - t_1$ , then the western at the day that a pared a, well therefore a mad B are the rate of a rationary parameter  $t \neq 0$ , and  $t \neq 0$ , and  $t \neq 0$ , are that  $t \neq 0$ , and therefore  $t \neq 0$  are that  $t \neq 0$ , and therefore  $t \neq 0$  are that  $t \neq 0$ , and therefore  $t \neq 0$  are that  $t \neq 0$ , and therefore  $t \neq 0$  are that  $t \neq 0$ , and therefore  $t \neq 0$  are that  $t \neq 0$  are that  $t \neq 0$ .

Corollary. Let  $a=\phi r$  t and  $P_1$   $\phi_1$  be the convergents of  $a_t$  then

$$t P_{2n}^{a} - \tau Q_{2n}^{b} = (-1)^{1+\epsilon} t$$
 (9, 11)

holds for every k = 1, 2, ...

I . Let s, be the compact fractions of a then a, - a, an

$$u_{xx} = u_{xx} + \frac{1}{u_{xx-1}} = 2x + \frac{1}{x_x} > x + \frac{1}{x_x}$$



But as 
$$a = \frac{P_{a+1}}{Q_{a+1}} + \frac{I_{a+1}}{Q_{a+1}}$$
,  $Ia = 13$   

$$a = \frac{P_{a+1}(a+a) + P_{a+1}}{Q_{a+1}(a+a) + Q_{a+1}}$$
 by is, hence

$$Q_{+a} = P_{+a} + P_{+a-1} + a(Q_{+a} + Q_{-a}) - I$$
, 0

 $\Lambda_{\theta} \rightarrow = \pi \div 1$  is rational, and a is irretional

$$P_{s+1} + P_{s+2} - Q_{sn} - r_{s-2}$$
,  
 $= P_{s+} + Q_{sn} + + Q_{sn-2} - n$ 

hold. If we multiply these equations by  $Q_n$  is respectively  $\rightarrow P_{nn}t$  and add, we get  $(P_{1n}^n \rightarrow r, Q_{1n}^n + t(P_{1n-1}, Q_{1n}), Q_{1n-1}, P_{1n}) = 0$  and from this formula we get (B, 11) directly.

#### 5 4. APPRICATIONS ON THEORY OF NUMBERS

It is proposed to selve

J 1,

[4,1]

by integral z and y

Of we and the country of the control of the control of and the control of the con

$$a = b = (x_1, \dots, x_{2n})$$

" bee Pan Ques and an a and b are positive and electrons prime

 $a = P_{\pm m} b = Q_{\pm m}$  and therefore

$$a \in Q_{2n-1} + b \cdot P_{2n-1} = P_2 \cdot Q_{2n-1} - Q_2 \cdot P_{2n-1} = (-1)^{-2n} - 1$$

ho de Hence we get the integral solutions by

$$x=Q_{2m-1}+k$$
 0, 
$$y=P_{2m-1}+k$$
 0 where  $k=0, \pm 1, \pm 2,...$ 

To solve  $x^2 + dy^3 = 1$  Pollow protein by integral which is well use 4/2 (3, 11).

$$q_1 d = q = q_1 q_1 \cdots q_n$$

then Pf -d Qf = (-1)\*\*

107,999.

22

ALUEBRA

Therefore if a to even.

$$\{\sigma_x,y\}=(P_{xx},Q_{xx})$$

and if w is old,

$$(x, y) = (P_{y+n}, Q_{y+n})$$

6

are so ut, by fice every positive its eggs to

E g. 
$$x^{0} - 20 y^{0} = 1$$
  
 $\sqrt{26} = (5, 10)$ 

Hy this method we got the solutions

$$(x, y) = (P_y, Q_y) = (51, 10)$$
  
 $\approx (P_y, Q_y) = (53013, 1020)$   
 $\approx (P_y, Q_y) \approx (530131, 103030)$ 

#### C. NTINLLO WILL THESE WILLIS FLEWINGS AND \$17.

[ 1] In that 4 the western of the positive annihers has been taken for the yelsem A. Schong the ring fitte integers. We will now consider another system A.

Let K) an arbitrary field and can indefinite n t included in K. The el mepts f is will be denoted by a b, f with and a thout is lower. The elements of the rong  $K > \frac{1}{2}$  will be denoted by

$$f(s), g(s),$$
 
$$(5, 1)$$

the ring to abining of a field) have well be not as the ring S.

In refer by t a new system A we create new elements denoted by threek

 $\phi(x), \psi(x), \chi(x), \omega(x)$  (5, 2)

in the following manner.

$$= (r_1 r_1 + r_2 - r_3) e^{r_1} + \dots + e_{n} + e_{n}$$

This is a pure a formal del a tion. It means that to every sequence of our thought from h with fixed decreasing integral indices

Say Camb sees



then g in a tole open of the new characters and the element will not be characterful withdress, a tole set tour because the Mining Landon representation about a fixing the We have to define the election and the multiposition of the electric t-2 and we will the term in such a way that the characters t-2 for weak t, we fix t < 0 form a submap isomorphic to K [x].

So we define a

Let 
$$n \ge m$$
,  $\phi(x) = \sum_{k=0}^{m} a_k x^k$  
$$\psi(x) = \sum_{k=0}^{m} b_k x^k = \sum_{k=0}^{m} b_k$$

then 
$$\phi_{-}(x) + \psi(x) = \chi(x) = \frac{\pi}{2} c_{\pm}x^{\pm}$$

$$\phi(x),\ \psi(x)=\omega(x)=\frac{1}{2}\mathbb{Z}_{+}^{2}d_{+}x^{4}$$

 $n \ge r \ge k - m$ .

$$a_1 = a_2 + b_1, \quad d_4 = X a_1 \cdot b_{n-1}$$
 (5, 4)

The detailions ( ) are breatly independ it of pull coefficients put before, the commutative, associative and distributive new held, and the subtraction is uniquely defined by

$$b_{\alpha} = a_{\alpha} - a_{\alpha}$$

Themse for the multiclement every constraints it. If for the constraints of  $\varphi(x)$  the constraints  $\varphi(x) = 0$ , for a 14 hard, then  $x \neq (x + y) = \sum_{i=1}^{n} x_i^{i}$ 

The elements r = 2 form a ring R and those elements for at h = r, r = 0 form a sutring for who have add to a and multiplication has been if include the same way as for polynomials. Hence there is an isomorphism l by which is a subring become asomorphic to h(x). Let  $\phi(x) = r$ , then  $\phi(x)$  has at least one co-sub-rank l 0, let a be the high stands of the periods l be co-efficients, then

$$\phi(z) := \sum_{i=1}^{n} \alpha_{ij} z^{ij} z_{ij} \oplus \Omega_{k}$$

non-man territhm egrowith , from 12 tf flowed rectly

The degree of a product + qual to the sum of the degrees of the feature.

The degree of a sum of elements I different degree is equal to the mostmum degree of the summands

be in the remark of the a finite product of two elements \$11 cannot be exact of these thereing of the consists of 2 man error. We will now more the coments of kind of with the entropy in large tenents, i. 2.

Bo the elements  $\sum_{i=1}^{\infty} e_i e^{i\phi_i} h_i$  to first otherwise dentified with  $h_i$ 

Line by  $(r-r^{-1}-r^{-1}-r^{-1}+r)$  then  $r^{2}+r^{-2}=1$ . Every field containing 1 into an input 1 to 1 the remarks of 1, which are  $r^{-1}$  etailed and  $r^{-1}$  is a form a ring which a none in the 10 mailtaining of Q. It is means of this ring will the form be identified with the responding consents of Q. So, r recommends in the field with  $r^{-1}$  and the responding  $r^{-1}$  is because lengthed with the second with

sum 
$$\phi(x) = \sum_{i=1}^{m} x^{i}$$
, where  $a_{i} = 0$ , for  $k < -m$ .

the polynomials to the elements (5, 2)

their highest en-efficients,

By his stion of this principle we get

$$e_+ \circ = e_+ \circ e_- = (-e_+ \circ e_-) \circ e_+ \circ e_- \circ e_+ \circ$$

and to further tope to make get an an engineerable set if comente co. IK

$$L \leq n + m_{\phi} \leq 1) \text{ or } - \chi_{\phi} \approx \sum_{x \in A_{\phi} \in B} x^{x} \text{ and }$$

$$\phi_{i} = x + \phi_{i} = \chi_{i}(x) \psi_{i}(x) - \chi_{j}(x) \phi_{i}(x) \operatorname{degree} \phi_{i+1} < \mu_{j}$$



then we crow references and the second statements of the second statements and the second statements are second statements and the second statements

For an inequality of K  $[\,\sigma\,]$ 

$$I \cap \phi \circ r = \sum_{i=1}^{n} \sigma_i \circ -\sum_{i=1}^{n} \sigma_i \circ r = f(\sigma) \circ \phi_i(\sigma) \qquad [5.3]$$

This supreme time of a consequent of permitted as a and an entirement which is a consequent of the con

And I we will be the setting that the property of the setting of

So we can apply a for 1 for 1 of the case changed dements

of V in the elements 5.2 follows:  $A = \{x_1, x_2, \dots, x_{k-1}, x_{k-1}\}$  are the polynomials in  $x \in [Sec (1, 1), (1, 2)]$ 

If we may the any ordered by the series of the property of the of the officer of the officer of the ordered by the series of the

The series the community of a management of a management of the series o

for degree of a polynomial has just a rate property on the forest on No. 10 ft. a from that the control by a fine the rate we got therefore the

and every element of Q a regrescrated by the great and great from the

We will now use the representation of coal numbers by continued fractions in order to it, is many the claiments in 2.

[8/8] As (1, 15) holds,

$$a_1 - \frac{P_1}{Q_1} = \frac{1}{a_0} \frac{a_{1-2}}{Q_2}$$

we will prove that the right a le of this citual in the an element. Whise degree de reuses in tefinitely is kincremes.

The lements is base then before the  $a_1$  is  $a_{i+1}$ , a common factor being arbitrary (see (1, 12)) hence  $a_0$  as  $a_{i+2} = a_1$  is therefore

$$\operatorname{degree}\left(x_1 - \frac{P_1}{Q_2}\right) = -\operatorname{degree} Q_0 - d_0 - \dots - \frac{1}{2}, \qquad (5, 5)$$

 $d_* > 0$  for k > 1 and  $a_* = x_* + 1 - a_{*+1}$  (see 1, k = 1). As the dispress of an of summan world fferent degrees in equal to the highest of the degrees of the normalism a degree  $(x_*) = \log (x_*) = 1$ , for k > 1.

$$Q_1 = 1$$
,  $Q_0 = s_0$ , ...,  $Q_r = s_r Q_{r-1} + Q_{r-0}$ .

Hence degree  $Q_{+} > 1$  for l > 1, hence the degrees normale with the ender and therefore

Hence from 3, 5 and 1 6 if f llows that

degree 
$$\left( \tau_1 - \frac{\mathbf{P}_{\tau_1}}{\mathbf{Q}_{\tau_1}} \right) - I = -2 - \sum_{i=1}^{n} d_i - I_{i-1} =$$

$$\operatorname{degree}\left(\frac{-1}{Q_{+}Q_{+++}}\right) \le 1 - 2 h. \tag{5.7}$$

The officency of the approx nation of an irrational number of it a continued fraction became country the theorem that if  $\frac{a}{b}$  approximates whether than  $\frac{P_{a}}{Q_{a}}$ , then  $b>Q_{a}$  is the corresponding theorem both in the case we consider here.



The could be and a some polynomia s of h [x]

$$H = \frac{f}{f} = \frac{P_{s,t}}{Q_{s,t}} \pm 0 \quad \text{and}$$

$$U = \text{degree} \left( c_1 - \frac{f(t)}{f(t)} \right) < d \quad \text{degree} \left( c_2 - \frac{P_{s,t}}{Q_{s,t}} \right), \text{ then}$$

$$\text{degree } f = 0 > f \text{ gree } Q_{s,t} \text{ holds}$$

$$(5.8)$$

Pro 
$$f = l = \begin{pmatrix} v_1 & P_1 \\ Q_1 & P_2 \end{pmatrix} \cdot \begin{pmatrix} r_1 & r_2 \\ r_2 & r_1 \end{pmatrix}$$
 The two summands on the right

nde have deferent degrees heure  $f = \log r + 3$  is  $\Theta_r = A_r + CQ_r + r + r$  polyz mosl,  $0 \le d$  gree  $x \in Q_r + r + r + d$  gree  $x \in Q_r + d$  gree  $x \in Q_r + r + d$  gree  $x \in Q_r + d$ 

 $\operatorname{degree} f \neq \operatorname{id} f + \operatorname{degree} Q_1 + d_2 f + \operatorname{degree} Q_2 f = \operatorname{degree} Q_2 f + \operatorname{degree} Q_3 f + \operatorname{degree} Q_4 f$ 

> degree Q, [use (5, 7), (5, 8)].

This theorem chaos a sa to of rox mate fund in a given by a power series of  $\frac{1}{x}$  by rational function in the neighbourhood of x=1 0.

Exercise. 
$$\frac{1}{2} \log \frac{x+1}{x-1} = x^{-1} \cdot \frac{1}{8} x^{-3} + \frac{1}{6} x^{-5} + \dots$$

Represent this function by a continued fract is and approximate it by rational functions.

Lemma. Let 
$$a = (s_1, ...)$$
,  $a' = (s'_1, ...)$  at m be the lowest index [5/4] for which  $s_m + s'_m$  holds and  $(s_1, ..., s_m = A, (s_1, ..., s_m = A'))$  then degree  $(a-a) = \text{degree}(b-A)$ , body (5, 0)

Prof. With it cas of generality we suppose that degree and degree and we shall use the indicate in data in for the convergence of a and for those of a' we shall use them with a dash.

$$P_i = P_i$$
  
 $Q_i = Q_i$  for  $i < m$   
 $Q_n = s_n Q_{n-1} + Q_{n-0}$  degree  $Q_n = r + q$   
 $Q_n = s_n Q_{n-1} + Q_{n-2}$  degree  $Q_n = r + q$ 

MIGEBRA

$$A = A' = \begin{pmatrix} 1' & & & 1 & & \\ Q & 1 & & & 1 \end{pmatrix} = \begin{pmatrix} 1' & & & + & \frac{1-1}{Q} & \\ & & 1 & & & Q & 1 \end{pmatrix} = \begin{pmatrix} -1 & & & & \\ & & & 1 & & & \\ & & & & Q & & Q \end{pmatrix}$$

$$= (-1)^m \stackrel{\pi_m}{=} \stackrel{\pi_m}{=} \stackrel{\pi_m}{=} .$$

1. If  $e = e^r$ , degree  $(s_n - s'_n) \ge 0$  $|s_n| > 45 - 5$ ,  $|s_n| = \frac{1}{5}$ 

2 If 
$$r > r'$$
, degree  $(r_m - r'_m) = r$ 

 $\Gamma_{\rm B} r = \chi - \chi_{\rm C}$  , r = 2 - 4 r (e.g.,  $r = \frac{1}{Q_{\rm en} \lambda}$ ). Hence

degree 
$$(A - A') \ge \deg \frac{1}{Q_{i,n}^{-1}}$$
 (8, 10)

#### holds in every case

From (5, 9) it follows that

$$d_{\rm k} n = 0. \qquad d_{\rm k} n + \frac{1}{Q_{\rm k} Q_{\rm k}} \sim \frac{1}{Q_{\rm k} R_{\rm k}} \sim \frac{1}{Q_{\rm k}} \sim$$

and that do so A  $\rightarrow$  1 sec  $\frac{1}{Q}$  Q = 1 such  $\frac{1}{Q}$  Hence

districts a degree [A-A A A A A A A A A A')

the degree of the first or of the live a name and is greater me the

and let t > 1 degree t > 1 is there exists a continued fraction of  $t_1 \neq t_2$ .

In the other and the summany the saids

$$P_1 = Q_1 \quad \text{and} \ P_{g_1} = (r_1, \dots, r_{g_r})$$

Let 1 < + < N

From the price log lessmant for we that here P Q Q = P a Q at

$$\equiv \operatorname{degree}\left(P_{n+1} \mid Q_{n-1} - P_{n-1}Q_{n}\right) \equiv \operatorname{degree}\left(\frac{1}{Q_{n}\mid Q_{n-1}}\right) \equiv -|\lambda_{n}|$$



where he increases to infinity, with a

$$P_{y} = Q_{y} = \sum_{k=0}^{\infty} a_{x} x^{k} \equiv \sum_{k=0}^{\infty} b_{x} x^{k} + \sum_{k=0}^{\infty} a_{x} x^{k}$$
 (5, 11)

The coefficients  $b_1$  are independent  $f > A_0 t_1$  increases with the index a we get an infinite set  $b_1 = b_1$  referring

$$\phi(r) = \sum_{k=0}^{\infty} b_k (r^k)$$

I mally we have to prove that  $\phi \neq i = \{i_1, i_2, \dots, i_{n-1}, i_n \} = \{i_1, i_2, \dots, i_n \} = \{i_1, \dots, i_n \}$ 

degree  $(\phi(x) = \{e_1, \dots, e_n\}) = \text{degree} (\phi(x) = \{e_1, \dots, e_n\}) \text{ bolds}.$ 

Find 
$$\phi(x) = (e_1, \dots, e_n) = b_{k_n} x^{-k_n} + e_{k_n+1} x^{-k_n+1} + \cdots + e_n = 0$$

di gros. A. and feere isns infini ety with a

# 6 CONTINUES PRACTICES WITH DATIONAL P. EMINTS

Let S be the field of the returnal principle then every finds a continued [6,1] fraction

$$\begin{aligned} (a_1, \, \dots, \, a_n) &= & \frac{P_n}{Q_n} = & \frac{P_1}{Q_1} + \begin{pmatrix} P_n - P_n \\ Q_n - Q_1 \end{pmatrix} + & + \begin{pmatrix} P_n - P_{n-1} \\ Q_n - Q_{n-1} \end{pmatrix} \\ &= a_1 + \frac{1}{Q_1} \frac{1}{Q_n} + c_0 + & \frac{1}{Q_{n-1}} \frac{c_n}{Q_n} \end{aligned}$$

represents a rational number, but an infinite continued fraction defines a number if and only if

$$2\frac{(-1)^n}{Q_{n+1}Q_n}$$
 (6, 2)

converges. If the sun ,6, 2) is convergent

defines a real number equal to (0.3). A necessary confirma for the convergence of (6, 8) is therefore

$$|Q_{n-1}Q_{-}| \longrightarrow 30$$
 [6, 4)

If the numbers  $Q_n$  are other > 1 on b, r < 0 such r is definiting and 1 on a the continued front in the right,  $Q_n = Q_{n-1}$ , increases at only and  $Q_n = 0$ , a satisfied

[6/2] Theorem. If Z | .. | is convergent, (6, 3) is divergent.

Proce. We prove by mathematical attret on that

$$Q_n < \prod_{j=1}^n (1 + \lfloor n_j \rfloor)$$
 (6, 5)

As  $Q_1 = 1$ ,  $Q_2 + e_0$  the fermula halfs for n < 3. If 0 - 5 is true for  $n < m_s$ 

$$\begin{aligned} Q_{n} &= s_{n}Q_{n-1} + Q_{n-2} \\ Q_{n} &\leq \prod_{j=1}^{n-2} (1 + j \cdot i) + (j \cdot i, j \cdot i) + (i \cdot i, j \cdot i) +$$

If  $\Delta \tau$  converges the  $\alpha \tau$  product H  $1 \tau$  [  $\tau$ , , conver,  $\alpha \tau$  to a post vector Q and [  $Q_{\tau}$ ]  $\sim Q$  is do for every order  $\sigma$ . Hence (0.4) due not have and the continued fraction a divergent

[6.3] Let s > 1 for > 1 for in Q<sub>1</sub> = 1 Q<sub>2</sub> = s<sub>2</sub> > 0, Q<sub>2</sub> = s<sub>3</sub> Q<sub>4-1</sub> + Q<sub>n-1</sub> (to see by mathematic son action that each number Q<sub>1</sub> > 0.

to 2, in therefore an alternating acres whose elements have atendity ducreus up at a late values. This series converges therefore if and only if it, 4 is matrix of. These mail texts us lead to the feet wing.

Theorem Let \*, > \* for \* > 1 then the automod fract on (6, 8) as convergent if and only if \$20, is divergent

Proof If \$1, = \$ a | is convergent the continued fraction in divergent, as it has been proved by the \$10 at no theorem.

<sup>&</sup>quot; We use the term " divergent. For every two convergent series



Let  $\Sigma_{x_i}$  be divergent then  $x_i \rightarrow \infty$  . As  $Q_i > 1$ ,  $Q_i = 1$ , and  $Q_{y_{i+1}} = x_{i+1} Q_{y_{i+1}} Q_{y_{i+1}}$  we the third of the  $Q_{y_{i+1}} > 1$ . Hence  $Q_{y_{i+2} + y_{i+1}} Q_{y_{i+1}} + Q_{y_{i+2} + y_{i+1}} Q_{y_{i+2}}$  and as  $Q_{y_{i+2}}$ , it follows by matter at that

$$Q_{2r} \leq \sum_{k=1}^{\infty} r_2$$

Hence

$$Q_{2n-1}, Q_{2n} \longrightarrow \infty$$
, and  $Q_{2n}, Q_{2n-1} \rightarrow \infty$ 

therefore 5.4) a satured and as we alsed above this condition a sufficient for the convergence of 6.3, in the case considered here.

Let spect specification that the continued from the present of the form the

$$\begin{pmatrix} P_1 & P_0 \\ Q_1 & Q_0 \end{pmatrix} = \langle 0, 1 \rangle$$

We will show that this value is rectional

then 
$$a_1 > c_0$$
 and  $\frac{c_2}{a_1} = \frac{1}{c_2 + c_3}$  hence  $c_0 = \frac{c_3}{a_3} + c_2 = \frac{a_3}{a_3}$ , where  $a_3 = a_1 + a_3 a_3$  is

integral as  $n_g = 0, s$ . ), there is  $< r_2 < 1$  bence  $a_g > n_s > 0$ . In the sound minimizer we get

$$a_9 = \frac{1}{4 + a_1}$$
,  $a_4 = -a_4 = \frac{a_4}{4}$  and  $a_5 > a_4 > 0$ 

By repet tion if this procedure we get an intime set of decreasing integral positive numbers

and that is impossible.

Hence of is irrational.

Bringh

$$x_1 = 0$$
  $x_2 = 2$   $x_3 = 1$  and for  $x_1 > 1$ 

$$\epsilon_{g} = \frac{2 \cdot 4 \cdot 2 \alpha + \frac{2}{3}}{1 \cdot 2 \alpha + \frac{4}{3}} > 1 \cdot \epsilon_{g + 1} - \frac{4 \cdot 7 \cdot 2 \alpha + 1}{2 \cdot 4 \cdot 2 \alpha + 2} > 1$$

then it is

$$Q_1 = 1, Q_2 = 2, Q_3 = 3$$

att life in the Hentit is

$$2 - 1 - 2m = 2 - 1 - 2m - 12 - \epsilon_{2m} + 2 - 4 - 2m - 2$$
 
$$1 - 3 - 2m + 1 = 2 - 4 - 2m - \epsilon_{2m+1} + 1 - 3 - 2m - 1 ,$$

it follows by mathematical resuction that

$$Q_{gm}=2$$
 4. .2m,  $Q_{gm+1}=1$  . 3,..2m + 1

b mee

the angletical fraction is treatment. Its visco a

$$\frac{1}{Q_1} \frac{1}{Q_2} \frac{1}{Q_3} \frac{1}{Q_4} \frac{1}{Q_4} + \frac{1}{Q_4} \frac{1}{Q_4} \frac{1}{Q_4} + \frac{1}{Q_4} \frac{1}{Q_4}$$

Hence e in irrational



# PART IV. APPROXIMATE SOLUTION

## Новмин'я веними

but k to an urbitrary field,  $\epsilon_{J}$ , the field of the real formers of  $\{1/1\}$  of the constant numbers q be an extremal f K and f(z) a position of K[x],

$$f(x) = \sum_{i=0}^{n} |\phi_i(x)| = (x-q) \sum_{i=1}^{n} |a_i(x)|^{n-1} + |a_i(x)| + |a_i(x)|^{n-1}$$

then

holds. We can arrang the care nation of the same tentes, a fill wa

We can tend it is sychosome method I, it will shop

$$f_1(x) = (x-q) \ f_{1,1}(x) + a^{\mu}_1 \ , f_{2,1}(x) = \frac{2}{x-2} \ (-x^{1-2})$$

After n-1 at the we get i x represented as a polynomial in x-q

the protection seknows in Analysis as the Taylor sent of the at the protection of the coefficients on cools be done on using the last me of 1 1' up to a part the compute scheme by the collection in called Horner's sub-me. It is the most convenient method

for calculating r > d q and the coefficients f > r are given and furnishes the representation of r by r = qt. The calculation will be expected by the following

Example:  $f(x) = x^4 - 35x^3 + 68x^3 - 119x + 67$ 

$$q = 1, 1 - 13 - 68 - 118 - 65$$

$$1 - 14 - 54 - 68$$

$$1 - 13 - 41 - 2$$

$$1 - 13 - 41 - 4$$

$$1 - 1 - 20$$

$$1 - 1 - 20$$

$$1 - 1 - 20$$

$$1 - 1 - 20$$

$$1 - 1 - 20$$

$$1 - 1 - 20$$

$$1 - 1 - 20$$

$$1 - 1 - 20$$

$$1 - 1 - 20$$

$$1 - 1 - 20$$

$$1 - 1 - 20$$

$$1 - 1 - 20$$

$$1 - 1 - 20$$

$$1 - 1 - 20$$

$$1 - 1 - 20$$

$$1 - 1 - 20$$

$$1 - 1 - 20$$

$$1 - 20$$

[1,2] If recent heights very small for calculating the roots. The method will be appared by the hold for a live assumption

7 1 L 7 L + 1 C 21 L + 1

$$x = q + 1$$
  $(-r) + z = -r^4 + 17 x^4 + 27 x^3 - 21 y + 2$   
 $f(1) = g(0) = 2$   
 $f'(1) = g'(0) = -24$ 

We therefore against the there is a roll of the more real. In the next without the selection of and the selection we have the property of the proposed in the first terms of the selection of the

<sup>\*</sup> The sage of means oppose motive and to " We apply been an electrology theorem on each mount force one which we be could be 2 dy



As f is the formula of a continuous therefore f(x) by

Hence the cost is approximately a part to the liberal provides the cost of the solid strains and the cost of the

Hotto q = 0 000 6-124 - 1 000 124

On using this hipe simulates we can be easily get some more figures of this decreed development. The most differ these is a constant a to get a first approximate in a the roots. For the quarties of the chain product the values of -x that a soluble set of values x. We may get these values by Hornes a scheme, but it is after applied to abbreviate this achieve in following manner.

Oreca q que que en Woodbunte la Hance a schettre [1],

$$f(x) = b_1 + (x - q_1) \ f_1(x) \qquad f(q_1) = b_1$$

$$f_{1-x} = b_2 + (x - q_2) \ f_{2}(x) \qquad f(q_1) = b_1$$

$$f_{m-1}(x) = b_m + (x - q_m) \ f_m(x)$$

$$f(x) = b_m + (x - q_m) \ f_m(x)$$

$$f(x) = b_m + (x - q_m) \ f_m(x)$$

$$f(x) = b_1 + (x - q_m) \ f_m(x)$$

$$f(x) = b_1 + (x - q_m) \ f_m(x)$$

$$f(x) = b_1 + (x - q_m) \ f_m(x)$$

$$f(x) = b_1 + (x - q_m) \ f_m(x)$$

$$f(x) = b_1 + (x - q_m) \ f_m(x)$$

$$f(x) = b_1 + (x - q_m) \ f_m(x)$$

$$f(x) = b_1 + (x - q_m) \ f_m(x)$$

$$f(x) = b_1 + (x - q_m) \ f_m(x)$$

$$f(x) = b_1 + (x - q_m) \ f_m(x)$$

$$f(x) = b_1 + (x - q_m) \ f_m(x)$$

$$f(x) = b_1 + (x - q_m) \ f_m(x)$$

$$f(x) = b_1 + (x - q_m) \ f_m(x)$$

$$f(x) = b_1 + (x - q_m) \ f_m(x)$$

$$f(x) = b_1 + (x - q_m) \ f_m(x)$$

We will feel pather and read to the previous rumphers this manner.

$$q_{0} = 2$$

$$q_{0} = 2$$

$$1 -15 -08 -119 -67$$

$$1 -14 -34 -05$$

$$2 -14 -17 -2$$

$$4 -15 -17 -2$$

$$4 -17 -17 -2$$

$$5 -17 -17 -3$$

$$t(t-2-x-1+s)x-1-x-1+x-1+x-2+8-x-17, \dots (t-2-t+2+8-x+17, \dots (t-2+8-x+17, \dots (t-2+8-$$

The equation of each tilder x > 0, x > 0, y = 0, the lat, the 4 h and the each tilder x > 0, y = 0, y = 0, the lat, the 4 h and the each of the two other two late the y > 0. So a contact dim tilder extend 1 by We can sate 1 hor root in the attribute 1 2 throw and level in rad in the attribute 3 throw and level in rad in the attribute 3 to y = 0. The resident may all affect the by Harnor's achemic on an expression

[1] If We will use here a different method of an collection. If for a  $\sum_{i=1}^{n} e_i e^i = e_i$ ,  $e^i = e_i$ 

If  $\xi$  is a root of g and  $h = \eta_1 < b + 1$  by repolition of the procedure we wanged a representation of  $\xi$  as a continued fraction. If we step the calculation after a steps we get the approximation  $\frac{1}{2}$ , and the error becomes

$$\xi \frac{P_*}{Q_*} < \frac{1}{*Q_{*+1}} < \frac{1}{Q_*}$$



## HORNER'S BCHEME

This method will be a untrated by the example previously used. We know that  $x^4 - 15x^3 + 68x^3 - 114x + 87$  has a rad in the interval  $\theta \approx -1$  becomes we represent this polynomial by Horner's method as a polynomial in x = 8.

Hence 117  $q_1^4 = 137 \ \eta_1^2 \Rightarrow 92 \ \eta_1^0 = 17 \ \eta_1 = 1 = 0$ 

By manufact heart we see that l = q + 2 the fast conformation are positive, but l = q + 1 there is an approximation of the interval  $1 = a_{11} + a_{12} + a_{13} + a_{14} + a_{14$ 

In the same manner as it has been done for you see that we is amounted in the interval (1, 2)

## ALÖFBBA

go to an the interval (2, 4)

Property the reason has very an experience that quantities be obsequed to the time and of the last of the region of the other policy of an extra the region of the sound of the other the sound the

The next q will become = 1

Hence  $\xi = (8, 1, 1, 9, 4, 1, ...)$ 

$$P_{1} = -8$$
  $Q_{1} = -1$ 
 $P_{2} = -9$   $Q_{2} = -3$ 
 $P_{3} = -9$   $Q_{4} = -3$ 
 $P_{4} = -43$   $Q_{4} = -5$ 
 $P_{5} = -189$   $Q_{7} = -22$ 
 $Q_{7} = -23$ 
 $Q_{7} = -27$ 
 $Q_{7} = -27$ 



As the standard of the standar

In the last subsection Horacz a scheme two can be I for the [1/5] a lution of equal, we will reast and instance by restricts. In this above on temperature of the West and the stantant of the ground it the properties of the West and the stantant of the stantant of the contract of the stantant of the st

Hence if a is a root,  $\hat{\Sigma} b_s = \hat{\Sigma} b_s a^{-1}$ 

Let b, be positive name or and complex

- 1. If he at the equation countries countries to
- 2) If | a | = 1, a = 000 P+coin P,

 $\Sigma b_s = \Sigma b_s$  completely the property  $L(s) + 1 = b_s > 1 = b_s$  and the constant  $L(s) + 1 = b_s > 1 = b_s$ . If therefore  $m(a_s) \circ m(a_s) \circ m($ 

$$a_* < a_4 < < t_*$$

Then for every rant a of the parameter  $y = \frac{1}{2}$ , the runts of 2a,  $e^{x}$  satisfy < 1. The theorem a known as

hetherpoon the remarks the contract about of  $\sum_{i=1}^{n} x^{i}$  have at about value < 1 if the coefficients satisfy the

Ü

AUGEBRA

#### 1 2. THE SOUTS OF REAL POLYNOMIALS

# [2/1] In this section

with routh tinfers with beamless real combers, a the same

let to a mit the numbers of our long to become the paste

H use a e a in real; a a in positive, . - a = co, and if s = 0, a

$$f(x) = a_0 + a_1x + \dots + a_{n-1}x^{n-1} + x^n$$
 . (2, 3)

can be represented by 
$$t x = \prod_{i=1}^{n} (x_i x_i)$$
 (4.4)

foce Part II [18 2]).

In Proof 1 at K be the first of each pointers, and -s be the roots of star I then K at -K at is the field of the complex quinters and there is an automorphism I fithe field overchamping swith a and leaving the resonant members upgatered for will not be altered by I, hence a well be transformed to a the theorem to true

and Prior If a sareal, area If a monotonic (x = 0.04) is a equal possible and a read as the in the first of the road numbers. As f(x) and f(x) have a continuous these polynomials have a continuous factor of particle degree. Hence f(x) is directly by f(x) and an therefore a root of f(x).

Corollary 1.

$$f(x) = (x + c_1)$$
  $f(x + c_2) = x + d_2 \cdot (x + b_2) = x + d_3 \cdot (x + d_4)$  (2.5)



the fire a finite of the second in this period is a second of the second

former, 3. If a is odd, there exists at least one real root

on only I the factories of of fr. As P. or a symmetric polyto multiple, and with integral conflictors of factories and last II [40]. Only

$$\Delta = g \cdot (a_1, \dots, a_n) \quad (2, 1)$$

where queryo mail with integral or florents. From 2, 6 ( f. was two Part II, [10/8]) that

$$\Delta = 0$$
 if and only if  $a_1 = a_2$  for  $i \neq j$ . (2.8)

Let the a roote i, be all I fferint. As A has been proved in Part II [1 4].

To get divertage to interchange every nizer that has been pate from (2 ) the conjugate from new to the hat get of the new the determinant (2 )

Hence  $\delta = (-1)^{4} \lambda$ , and therefore

$$(-1)^4 \Delta = (-1)^4 A^6 = A^4 > 0$$

Hence the following theorem holds.

There is fit to dispute a base a filteren resta. Then he describe note of parties of par

findings. A res princement of degree 3 has three real roots don't only if the discriminant is positive

A real polan small of digree 4 with positive discriminant has either for rightfor nt real roots or two pairs of conjugate complex roots.

have see Prove the preceding theorem without the he p of 2 8

44

#### ALGEBRA.

- [2,2] If we say the notes that the property of the state of the state
  - 1. If a < b and the  $a_a a = b$  by a + b by different then there is a root of flet in the interval (a, b)
  - 2 H + c b and c + c + c + c + b so at of the interval (a, b)
    - 3. If  $f(x) = (x-a)^{k} g(x)$ ,  $g(a) \neq 0$ , k > 0, then

$$f'(x) = (x - a)^{x-1} g_1(x), g_1(a) \neq 0$$

for, g, legiere (o-al g'(s)

There is the second tens of a second interest tens of a second interest into the second interest

The property half (r v ex non y c to p n with a fin non-er of r t according to a property of y x none of the coefficients of the set of positive to, n we for very jet we very of r. Hen. f(a) a self-gase for a segmentary parameters product of \$ 7 form

organical party of a different problems of the signs in the sign in th

log decoping to sea ynomial a r h we get

$$f(x) = f_x(x-b) = a_{b+a} \cdot (x-b)^a + \dots + a_{b+a}$$
 (2, 10)

such that a second endependent of b 2 11



Biner a set if and real the parent set of the and from the w

$$\gamma_{b+a} = \frac{1}{b^{\frac{1}{4}}} f^{\pm}(b),$$
 (2, 12)

The sequence

tens ex reduced by making at overselement of the or II in the reasonable question, the amount of the form of the product of the normal of the form of

If the first and the not expent of terrelation expent boxs it is grouply a bit in expention of the first basis of the expense (i.e., b) is the call of the expense of the

If we strok not to 2 of an enquent which particles a hang be an different from the test be necessarily to make a test particle and the enquent of the particle and the enquent for a with a set if there is not to will compute the changes of a rows.

whore

Let q > 0. If a continuation a change in the 2' row har 2 to and it is of the same man as a little of the same man as a little of the same has a little true which have the same and the formed by these continuation of the second real lite has been as the creep poling elements of the second real lite number of changes of the 2' row is those for less or equal to the same of the changes on the first. The same had a for every pair of a long or town a the following system which we get by Historians between

But the 'ast row is identical with

Homes for q > 0 C 1 5 C h to

(2 - 15)

Lot c=b+q>b

$$f(a) = f_a(a - b) = f_a([a - b] - q)$$

then the number of charges in  $f_{ab}$  that the number of charges in  $f_{ab}$  is not greater than the number of charges in  $f_{ab}$  the

$$C(c) \le C(b)$$
 (2.10)

had a good a have the same number of changes we charges a mith of any one I generally touch 1, >0. If it was a rest une of a repention of a second of

## Throrem C decreases steadily to 0

for  $b < \min$  at these values be not expected by a root for then  $a_1 = f(b)$  and  $a_2 = f(b)$  have the cause again there  $C = b + 1 + 2k_1$  where  $k \ge 0$  is an integral number

to therefore a decent mount decreasing function taking only integral taking and the safety at pastes which are all rolls bave even walness. We have now to investigate the entire in the costs of fire.

Let c be a root of multiplicity we, end c-c-y, then

Lety Lac France

$$f(x) = f_{a(x)} + f_{a(x)}^{-1} + \dots + g_{a(x)}^{-1} + \dots + g_{a(x$$

Let  $\phi_{x,m}y^m = \phi_{x,m} (x + \phi)^m = g_1(x)$ 

If , >0 there to no change in the coefficients if g !!

If s < u, the electric cuts here alternating a gas and have therefore in changes. If r is small enough the last m+1 cutile cuts of  $f_{s}(u)$  have the same u gas as g > 0. If therefore u normals from a nell negative to small positive values, the number of the changes in  $\sigma_{d,u} = \sigma_{d,u}$  do reuses by an



even value and the number of the changes house of a discussion by m.

The he put there are not real name get the fill was bregget a trace

Theorem Is fide 100 to \$ , been ed by the Combine of the coats of form the piercal 0 and the resulting theorem multiplicity, then

$$C_1(h) = C_1(e) + e + 2k$$
,  $\bullet$  (2.18)

where k≥0 is an integral number

Apply  $r_{\rm R}$  the horizon of it constructs when a court  $r_{\rm R}$  great that  $C(r) \approx 0$ , we get an corollary

Iterative rule. The number of the positive roots of / configuration to the number of the positive roots of the number of the configuration of a run tracket from the number of supplies from the numbe

if we consider that in (2, 12)  $a_{evp}$  and  $f^{**}$  is have the same eight to 18 can be expressed in the first anglements.

the form F warrach area Lity h \$10 f - \$ 5. then the newtre of the contact f a in the interest has been great being contribute to win maltiplicity or equal to the other of the changes forms in the sets.

if control is the companies and and

If we set  $x = \frac{x + b}{b + 1}$  and therefore  $y = \frac{x - b}{-x}$  then the positive rects of

There free ulse I as a water give derically the exact puncture is be restricted an aternal, but they are very asciulate goving it exclude there complicated cases.

We will go total orther fetal and the assert e considered in 1

$$f(x) = x^4 - 15x^2 + 68x^2 - 119x + 67$$

#### ALGEBRA

As we stablished the rear roots in postive and so not do not be true of the being a closure of the best of the being a closure of the best of the best of the best of the contract of the best of the contract of the area of the best of the contract of the green of the best of the contract of the green of the best of the contract of the green of the contract of the signs only—

We know they rests, in in the interval 1.2 quistion in the autereal 8.5 ments of the case to the rest of particle and are the charge. We should find that the case the changes on a 8 corresponds to make of an interval particle was systematical these expected roots by Herrica a become and go be very a myle calculation.

$$C(3)=1$$
  $C(2.6)=8$   $C(2.6)=8$ 

Hence the two research is y be aduated in the interval 204 < x < 215 but we were proved that for the matters in the interval. We stated providedly that

Here a because the two racts and also be The same result an attachment of call delengt a case or same and stating that it is negative

[2/8] A the been given by Starr We suppose that for the local field by the been given by Starr We suppose that for the local field by the suppose that for the local field by the suppose that for the local field by the suppose that the suppose that the suppose that the suppose that the suppose the suppose that the suppose the suppose that the suppose that the suppose that the suppose that the suppose the suppose that the suppose the suppose that the suppose that the suppose

The method uses a come Store a chain of polyn in an

$$f\left(x\right) = f_{+}\left(x\right), f_{\pm}\left(x\right), \dots, f_{m}\left(x\right),$$

 $\{2, 19\}$ 



and the number C o of the changes of a grade

The chain should be a few such a manner that C = b is a easily decreasing function changing a value only at the rate of f(x) and having it there points a satisfable value. Then

becomes the homber of the robs of the anche of terral hi

For this purpose we have to arrach the chain u = t at at each r = t of f cone change with that and that in the control f > 1 the man we of the change u = t and u = t. In or u = t of the change u = t to u = t. In or u = t of t > 1, t = t of

(1) for (x) has the same sign, as f'(x)

In refer to as dominater as sometimes of the formation of  $f_{II} = f_{II}$ 

f<sub>m</sub> (z) should have constant sign, and

But a will have the sound country mout demogrations before and after passing a root. The country is it is seen a subsequence a this following theorem:

the number of the costs of the many interest by the second of the costs of the many interest by the second of the costs of

In order to get a line of the kind we may use the abjustifier to f

$$f_{ij}(x) = f^{ij}(x), \ f_{ij}(x) = g_{ij}(x) \ f_{ij}(x) = f_{ij+0}(x).$$

where he preced fith privation of the contribution

The fraction continues are organized as the compact factor of  $t_i$  as and  $t_{i-1}$  at an interval of  $t_{i-1}$  and  $t_{i-1}$  at the continues factor of  $t_{i-1}$  and  $t_{i-1}$  at the continues factor of  $t_{i-1}$  at the continues factor of  $t_{i-1}$  and  $t_{i-1}$  at the continues factor of  $t_{i-1}$  and  $t_{i-1}$  at  $t_{i-1}$  and  $t_{i-1}$  are the continues factor of  $t_{i-1}$  and  $t_{i-1}$  at  $t_{i-1}$  and  $t_{i-1}$  are the continues factor of  $t_{i-1}$  and  $t_{i-1}$  and  $t_{i-1}$  are the factor of  $t_{i-1}$  and  $t_{i-1}$  and  $t_{i-1}$  are the factor of  $t_{i-1}$  and  $t_{i-1}$  and  $t_{i-1}$  are the factor of  $t_{i-1}$  and  $t_{i-1}$  and  $t_{i-1}$  are the factor of  $t_{i-1}$  and  $t_{i-1}$  and  $t_{i-1}$  are the factor of  $t_{i-1}$  and  $t_{i-1}$  and  $t_{i-1}$  are the factor of  $t_{i-1}$  are the factor of  $t_{i-1}$  and  $t_{i-1}$  are the factor of  $t_{i-1}$  and  $t_{i-1}$  are the factor of  $t_{i-1}$  and  $t_{i-1}$  are the factor of  $t_{i-1}$  and  $t_{i-1}$  are the factor of  $t_{i-1}$  ar

In Street a method to a new avaposance to get the exact number of the forest roots to any interval out the peace and call and it constitutes were burieful and it is from the convenient to use Budan Fourier's to in its affect a with special considerations as we direct in the presence of a superstance to same example with Street method as an exercise.

Remark Starm's chain (2, 10) can a ways be replaced by

[274] Stormer the seem were with a patential I good a polynomials

$$P_{\alpha}(x) = \frac{1}{\pi} \left[ -\frac{1}{2} \left( x^{2} - x^{2} - x^{2} - x^{2} \right) \right] + \min \{0, 1, 2, \dots \}$$
 (2, 21)

D" fen talle a der, of fill and motten a ] and D' to the function stout | If a and man polynomers are

$$\chi(r_{(0,0)}) = \sum_{i=1}^{n} -(1)^{i}$$
 (2, 22)

where T m (2, 22) it follows for m>1

$$1 - \left[ -\frac{3}{2} - \frac{1}{2} \right] + 2m \times \left[ \frac{m-1}{2} - \frac{m-1}{2} + \frac{1}{2} - m + 1 - \frac{1}{2} - \frac{1}{2} + 1 \right] = \left[ -\frac{3}{2} - \frac{1}{2} - \frac{$$

O the ther hand we get from 2 22 for mod 1

$$2[(x_1 - 2 - y_1)^{-1} - 2y_1^{-1}] + 2[(x_1 - y_1)^{-1} - 2y_1^{-1} - y_1^{-1}] + 3[(x_1 - y_1)^{-1} - y_1^{-1}] + 3[(x_1 - y_1)^{-1}] + 3[(x_1 - y_1)^{-1} - y_1^{-1}] + 3[(x_1 - y_1)^{-1}] + 3[(x_1 - y_1)^$$

By softeach in A. 2, 25 from 2, 26 and applying 2, 21) we get

$$m(\Gamma_m) = e^{\frac{\pi}{4}} - 1 \| \Gamma_m - x \| + e^{\frac{\pi}{4}} r \Gamma_{m-1}(x)$$
 (2...2.)

2. 25 he do a so for re=1 and therefore generally

From

$$10^{m+1}\{(x^2-1)^m\} \approx D^m\{2mx (x^2-1)^{m+1}\}$$

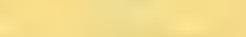
$$2mx(r^{-1}x^{2}-1-1)+2m^{2}(1)+1-x^{2}-1=1$$

we got

$$P_{m,1}(x) = x P_{m-1}(x) + m P_{m-1}(x),$$
 (2, 26)

brain a, 20 and 2, 21 we enternal I'm and go o

$$(x^0 - 1)P_m(x) = mxP_m(x) + m P_{m-1}(x),$$
 (2, 2t)



On replacing to by to + 1 in 2, 25) we get

4

$$m + 1$$
  $\chi(x) = (x^2 + 1)\Gamma_m^0(x) + (m + 1)x\Gamma_m(x),$  (2. 25%)

From (2, 27) and [ ] . . . . ) we choose a Pass) and we get

$$\tan + 4_1 P_{m+1}(x) = (2m + 1)x P_m(x) + \sin P_m - (x),$$
 (2. 28)

We consider the sequence

$$P_n(x) \cdot P_{n-1}(x), \dots, P_n(x) \cdot P_n(x) = 1$$

on the interval

$$-1 \le z \le \pm 1$$

From 2 =7 of f., as for mean that f.P<sub>n</sub> x at P<sub>n</sub> x booth a region  $P_n$  , x = 1 (P<sub>n</sub> x and f.x and five a common x = 1) be religious for a post of  $P_{n+1}x = x = x$  and there are form a 2 $\sigma$  and therefore of all adequate post on x = f. 2f in the state of the Henry top  $P_n(x) = 0$ ,  $P_{n+1}(x) = 0$ , and  $P_n(x) = 0$  of a overy work of  $P_n(x)$ .

 $A=B_{\rm m}(z)\sin(1z)$  ,  $r=2\cos(z)$  , which is forest the superstance of the superstance of

From (2, 27, it follows that

 $P_{m}(1) = P_{m+1}(1)$ 

In In Part 1

As Parze = 1, it follows that

 $P_m(1) = 1$ 

and

and

$$P_m = 1$$
  $-1.2$ 

The tour of langes of sign 2.21 is the time to 1 or 0.1 of themes, which is a district constant to be expected as a large free point many are all activities as the concern. Hence there are of large free points in a part all activities as the concern. — 1, \* 1 and are not persons.

To find out as stomat early the real costs of a polyto-most

[[ ]

$$f(x) = a_0 + a_1 x + ... + a_n x^n$$

with real recliments we have at first to find an otterval clotten and all these roots

Let

As that is for H — f is the same sign as  $\sigma_n$  on therefore a hand the same sign as  $\sigma_n$  of H and the roots of the roots of the intention of the intention  $\{-t,+t\}$ 

When it intervals he questi the real part of the parameter has been food as a major of the state and bed his to my rotate entire to the state of the

[2,0] The maked to a set the party of section to used by a charter some Pyth prote of section, ho was a teamer n what make terms make the first hours how there is no steethe polynomerous only for the entire like of the contract of the interest of the norm range containing only one margin rate the approximation on try to a set very quickly by Hornor's achieves.

function for many be replaced by the straight line connecting the party with observers a and b. Thus any intersects the country of the approximation of the control of the approximation of the control o

$$x - 1 - y - f = y = y^4 - 14y^3 - 20^{-2} - 24y + 2$$
  
$$g(0) = 2, g(1) = -3.$$

Hence we get by the real firm as approximation by

$$247^{2} + 94 + 197^{2} + 2997^{2} + 2 = 0.00$$
 (2.30)

H-nes



The approximates in tetter at the not good. As we call at the

$$x = 1 \pm 0.0935324$$

Of course the graph of the Polynomial is in that interval very different from a straight line. Now

$$g \circ = 2 \qquad j \in 4$$

$$g^{j}(0) = 24 \qquad j$$

The graph as there is a second to the state that the state and as remained to be required to be point so 1 = 0

The fermion to the descent in Horner execution have to be a stad. The method are a decreased in the destruction of the destruct

As for 
$$y_1 = 0$$
,  $y_2 > 0$  (82)  
and for  $y_1 = 0.1$ ,  $y_2 < 0.005$  bold.

and as you a contain us francis to for a line and the uniterest of the man was to which you and that is a rest. So New yor method a very useful on contain reach, but if the interval is big or the tangent makes on you amad angle with the axis. One mothed exemp by used

In the section of metimes reference has been made to the graph of a polynomial. So the reader may ask if just had without may not be leaded to get the rots of a polynomial. Of course insteads of the kind exist and are very helpfus to get a convenient first approximate a for the rots but in applying these over sole only these readers may succeed who are furnished with the theory and practice of mathematics as drawing.

and consecutive ructs of the three for bus a constant again in he nterval

<sup>• 4</sup> very union graph a method a confine exectant car method. Exercerence our Busherbach banez. Vor samogen a ser Agebre 19 144-14 t.

the latthes again to a them a a resing use attend by board there is a constant to a them are a second or a constant to a constan

If the same the following the contract of the administration of the contract of the contract

Rain  $a_{2}$  who main rate 0, explored the main terms  $a_{1}$  and terms  $a_{2}$  and  $a_{3}$  are some transfer at the first twin will enter the example  $a_{3}$  and  $a_{4}$  and  $a_{5}$  and  $a_{5}$  are the first transfer at the main terms  $a_{5}$  and  $a_{5}$  are the first transfer at the example  $a_{5}$  and  $a_{5}$  are the example  $a_{5}$  and  $a_{5}$  a

the southern to the south of the south the line to the with the south the line to the south the line to the south the line time to the with the south the line time.

\$2 is Lemma a suspicious nor of the fift wing the right

From a the matter  $h = e_{+}e^{-}$  for  $h_{+}$  tempty normal with restricted with respect to the part of  $h_{+}$  and  $h_{+}$  to and of the  $h_{+}$  applying and with respect to the physical with respect to the  $h_{+}$ 

$$g(x) = b_0 f(x) + b_1 f'(x) + ... + b_n f^{(n)}(x)$$

has not the afficent reactions to and the constraining proposition body when each real last took in mate3 with its as a maltiple ty

Here e for m=1 and m=1 to the formula let p> and it the the rem to proved for m=p. We safe that for m=p+1. If maked of h is then a  $\{0, \dots, p+1\}$  if maked of h is then a  $\{0, \dots, p+1\}$  if m is a  $\{0, \dots, p+1\}$ .



Let  $h_{ij} = x^{ij} + \cdots + y^{ij} + y^$ 

If we replace in his to be a trace of the presence of the graves pend ing derivates of feet, we get therefore

$$j(s) = g_{+}(s) - a_{-}g_{-}.$$

From the processing of the control of the control of the first of the control of

#### 1.8 GRARPER'S METHOD

By Gracepe to the dualities and the manufaction of separate [8.1] memoring at the manufaction of separate property quests which the manufaction of a second second the second second second the second second

Let 
$$b_1 > b_2 > b_4$$
 (8, 1)

be the ruote of the polynomial  $a_0 x^n + a_1 x^{n-1} + \cdots + a_n$  , then

$$\begin{aligned} & \frac{-i}{a_i} = b_1 + \frac{b_2}{b_1} & + \frac{b}{b_1} & b_1 - c_1 \\ & \frac{n_d}{a_1} = \frac{-a_n - a_n}{a_n - a_n} \\ & \frac{-i}{a_1} \frac{t_2}{b} \left( 1 + \frac{b}{b} + \dots + \frac{t_n}{b_n} + \frac{t_n}{b} + \dots + \frac{h_n - h_n t_n}{2} + \dots + \frac{h_{n-1}}{t_{n-1}} \right) \\ & \frac{-i}{a_1} \frac{t_2}{a_1} \left( 1 + \frac{b}{b} + \dots + \frac{t_n}{b_n} + \frac{t_n}{b} + \dots + \frac{h_n - h_n t_n}{2} + \dots + \frac{h_{n-1}}{t_{n-1}} \right) \\ & \frac{-i}{a_1} \frac{t_2}{a_1} \left( 1 + \frac{b}{b} + \dots + \frac{t_n}{b_n} + \dots + \frac{h_n}{2} + \dots + \frac{h_n - h_n t_n}{2} \right) \end{aligned}$$

$$\frac{\pi a_n}{a_{n+1}} = \left( \begin{array}{ccc} 1 & a_n & \left( \begin{array}{ccc} r_j & a_j & & \\ r_n & a_j & & \\ \end{array} \right)$$

$$\frac{I_{n} - I_{n}}{I_{n-1} - I_{n-1} - I_{n-1} - I_{n-1}} = b_{n}(1 + \dots)$$

16 6 A partity great for a 1 in the annihiter, one be something by the management of the approximation

$$F_{+} \wedge \frac{1}{n} \quad \text{for } i = 1, \dots, n. \tag{B. 2}$$

the right 2 is not a residence of the order and the product of the

Let  $x^n + a_1 x^{n-1} + \cdots + a_n = -a_n n$  has a the roote  $b = a_1 x^n + a_2 x^n + a_3 x + a_4 x + a_4 x + a_5 x + a$ 

tree to mente the first term of a between the property the specific terms of the specific terms of the specific terms of the property that the property of the

It, east tients of fact and about the tracker using achome

(1) 
$$a_0 = a_1 = a_2 = a_3 = a_4$$

$$a_0 = a_1 = a_2 = a_3 = a_4$$

$$(2) = a_0^{\pm} = a_1^{\pm} = a_2^{\pm} = a_3 = a_4$$

$$+ 2a_1a_2 = a_2a_4$$

$$+ 2a_1a_4 = a_4$$



# ORASPER'S METHOD

As in the first pair of himse corresponding a strategy differ only by the sign to indicate differ only by the signs in the 21 on. The purphers apprease very quickly therefore the conversable to make the fact figures denoting the decimals very contry. It is the purphers we see the use the notation

8\*456181 for 8:456181 . 10\* .

To extract the roots at the and of the calculation we need a gardhess. It is therefore used as to calculate more decimals than the tables of logarithms contain.

Entre ple,  $x^5 - 10x^3 + 16x - 2 = 0$ (11) 1 -1'0116 2 2556 -150 0 82 0.4 (2) 618 2230 "4"0234 440656 In 0 492 - 0.0344 (4) 48199 446112 16 - 1175729 9"12031 256 D 00032 + 0 00018 -17747979\*12616 (8) 1 10.00

In the next step the coefficients will become the squares of the preceding coefficients, and in no case the error will have officines on the first 'figures. Therefore we stop the procedure and on unto non the roots by the help of logarithms.

log, of the coefficients	log. r9	log ( a	2
0	7:24254	0.90332	8:0412
7 242 4	2 CHSDG	0.20063	1 6028
a 12.64 2.40224	1 09/049	0.13-08-1	0.1365
		0.30108	1050000
		= log 2	

The sign of the roots cann the determined by termific a method, we have to arrange a sportal acceptant of rather again in every case. In this example the coests onto have alternated, eight between the coests onto of f(-c) are therefore negative hence the roots of f(-c) are all positive.

For he king we from the elementary symmetric functions of the approximate roots, and we get

$$s_4 = 10, \quad s_2 = 15\,9999 \quad s_3 = 2$$
 for 10 16 2

[5.2] If a ceal polynomia beacomposite is to not the name a ways conjugate, and these have therefore the same also into value. Gracific a method has then are to be middle don this case. In example, will give valuable hints for necessary modifications.

We know from \$ 1 that the polynomial has two real roots  $b_1 \approx 7.5926$  and  $t_2 = 0.00324$  and two complex roots

The en culation by Graeffe's method is given " on the next page

If the procedure be a peaked, the two first and the two last coefficients will be any the squares of the corresponding coefficients of the line (6) that the third coefficient will dipend also up to the second and the fourth. We cannot expect that further repetit in of the procedure will make the third coefficient and expectate of his prophetors as two coefficient of the polynomial base an equal absolute value. If by a greater than the absolute value of

<sup>.</sup> An the sign in the 2' line is a ways we we notic these moss for abbreviation



# GRAEFFR'S METROD

the complex roots, then

$$\frac{-m_1}{d_{10}} = h_1 = \frac{(1+2h_1^m+h_1^m)}{h_1^m} \quad \text{or a suitable } m \ .$$

A rough mental calculation above that  $h_1^{\pm} \cup 160, h_1^{\pm} \cup 1^{\pm}2$  .

The same consideration made for  $f, \frac{1}{x}$ , shows that if  $h_4 < h_2$ 

Let B and I be two intervals within every number of C is very small a companion to the numbers of B and let

$$f(x) = a_{p_1} + g_3 \times + \cdots + a_0 \times n + r + s$$

a ,

8 mold e coefficients in order to get the law of dependence we sha generation the considerations.

I

have two sets of roots

let ex to then the her ance approximately equal to the men

symmetric fun ismental function of by b and

$$= \frac{d_{\theta}}{d_{\theta}} \cdot \left( \frac{d_{\theta}}{d_{\theta}} \right) \cdot e_{\theta}^{-1},$$

where the a suitably chase mean value of the roots of the second set, and sherefore a number of C. Let y be a number of B. then

tek —— I made to one of the costs by theo i, by the roots of i (r) to the roots of i (r) but as a common in the cost to coefficients to be approximated by the roots of

$$\lim_{n\to\infty}\frac{1}{n}\lim_{n\to\infty}\frac{1}{n}\lim_{n\to\infty}\frac{1}{n}=\frac{1}{n}\lim_{n\to\infty}\frac{1}{n}=\frac{1}{n}$$
 and 
$$\lim_{n\to\infty}\frac{1}{n}\lim_{n\to\infty}\frac{1}{n}=\frac{1}{n}\lim_{n\to\infty}\frac{1}{n}$$

the absolute values of h belong to an interval H, the absolute values of e', being to C and every number of e' is small in comparison to H. Hence the roots h, we be approximated by the roots of  $e_{n}$ , e' e' e'  $e_{n}$ , and therefore the roots h of f(x) can be approximated by the roots of

Bo the polto mink x has to be apart up in two polynomials the first in defined by the r = 1 upper terms and leads to the upper class of roots,



# GRARFER'S METHOD

the second one is defined by the silver terms, and sads to the lawer class of roots.

The two classes may also be invided into sob-seases ste. Finally we get discuson

on he fact being small in a migration with the costs of the preceding classes, and to each class electropic is a point of all which can be out out from f(x). The rates of the absolute value of the roots in reason when we replace these roots by higher powers of them therefore we get finally by Granife a most of k polynomics as a bird them having only for is with the notice absolute value. In the previous cosmoles these polynomials are

$$x = 1716371 - 1 - 171 x^2 - 279413 x + 4738747, 4779767 x = 246$$

From these polynomias we get the code of the

$$= 0.79767 \qquad \log_3 h_0 \left[ -0.22466 - h_0 \right] + 1.6776$$
 
$$\log_3 h_0 \left[ h_0 - \log_2 h_0 \right] 10071 - 9.43265 + 10$$
 
$$8p = +79781^p + k \cdot 800^p$$
 
$$\phi = \pm 9711928^p + k \cdot 457$$

To fin shifthe calculation we have to fix the eigen of the real costs and to determine the entegral number k. As the eigen of the coefficients are alternating there is no negative test. Hence  $h_1 = 7$ ,  $h_2 = 0$ , the electric problem is a support of the results of annual m = 1. However method and by Lagrange's method

$$h_0 = 1.0502 + 10.50000$$

# AT GEBILA

For checking 
$$b_{1} + a_{2} b_{4} + 2 c_{3} r = 0.50104$$
  
for  $\log 2 = 0.80108$   
 $b_{1} + b_{2} + b_{3} + b_{4} = 100000$ 

If we tep a in the last work along the value for by by the more exact value distributed or more 1 by - 7 1200 we will get

 $h_3 + h_2 + h_3 + h_4 + h_6$  In sense the result can be corrected by firther a matern As we are from the results of I and from the characters, given here  $h_1, h_4$  and rare very exact. The correction is there fore experted to a near mainly the angle a whose true value may be a lattle and er. As a twent is a small angle this correction will mater  $ah_2$  and b are true up to the second decimal only.

If a p years a with roots of equal absolute value has a degree > 2, rather than multiple roots or at has a mach, again roots. The multiple roots who be some ved, when we divide by the first of the polynomial and its directive. You compagate ratio of equal absolute value can be elemed away by Historian scheme, its if |x| = |x| and x is different from x and x, then  $|x-a| \neq |x|-a|$ .

Hence the real and the compact mots of the can be found not by a combination of treaches method and Horner's scheme in every case. The results should be varified and it is passed to unminuse the error by the methods given in § 1.

## 4. ROOTS OF COMPLEX POLYNOMIALS.

Let 9 2) he a polynomial with complete cor becota,

$$\phi_1(x) = a_0 + a_1 x + \cdots + a_n x^n$$

$$\phi_1(x) = a_0 + a_1 x + \cdots + a_n x^n,$$

$$\phi_1(x), \phi_2(x) = f_1(x), \quad \phi_2(x) = f_1(x) \phi_1(x)$$

$$\phi_1(x) \phi_1(x) = f_2(x),$$

then the undity a new rest polymers als. On applying Graeffe's incitted to these polymers als we get the roots of \$20, but out of two conjugate



roots of  $f_0(x)$  have a root of  $\phi(x)$  and we have therefore to make a verification finally

Let 
$$||x|| \ge \sum_{i=0}^{n} ||\frac{y_i}{x_i}|| = i$$
, then

$$\frac{1}{a_n} \left[ q(s) \geq -z^{\frac{n}{2}} , -\frac{\sum_{i=1,\dots,n_n}^{k-1} a_i}{z^{\frac{n}{2}} + \frac{n}{2}} \right] z^{\frac{n}{2}} - z^{\frac{n-1}{2}} \quad (-n-1) = n^{\frac{n-1}{2}} > 0$$

hence  $\phi$  (i) \$ 0, and then fore the absolute value of the roots of  $\phi$  z is < f.

Another I is therefore the roots can be found out by hiskeys a theorem,

The motor of p is an entermode of the real polynomial  $l_1$  is  $l_2$  if q in f(x) where x = x + a and the real number a can be chosen in each a manner that  $f(x) = a_1x^2 + \cdots + a_1x + a_2$  has positive cost cost and  $f(x) = a_1x^2 + \cdots + a_nx^n + \cdots + a_nx$ 

For the polynomia s with positive coefficients the for wing theorem holds

Theorem Last the confidents of  $t = a_1 + a_1 x + \dots + a_n x^n$ be positive and  $0 for <math>k = 1, \dots, n$ , then the roots of  $t \neq 0$  have to satisfy the nonlinear

Proof Let 
$$x = q u$$
  $h(x) = j u = \sum b_i u^i$  then  $b_i = q^i u_i$ 

Hence  $b_{x+1} - b_x < 1$  from hakeya's theorem it follows therefore, for the roots that |y| < 1, and |x| < j

The roots of  $F(z) = u_0 e^z + ... + u_{n-1} z + e_n$  are reciprocal to the roots of i at As  $\frac{i_{n-1}}{u_{n-2}} < \frac{1}{p}$  bids, it follows from the first part

of the proof that the couts of 8 have to extisfy

$$|x| < \frac{1}{p}$$
 Hence  $|x| = |\frac{1}{x}| > p$  holds for the roots of  $|x|$ 

An interesting connection between the roots of  $\phi = and$  is derivate  $\phi'$  (a) is given by the

The sem : there Every convex polygon menuding all the roots of  $\varphi(x)$  contains every root of  $\varphi'(x)$ 

Proof Without any loss of generality we can supplie that  $\phi$  and  $\phi$  have no contained that  $\gamma$  be an arbitrary riot of  $\phi$  and  $\beta_1, \dots, \beta_n$  be the roots of  $\phi$ , then

$$\frac{\varphi^{-1/2}}{\varphi^{-1/2}} = \frac{1}{2} \frac{1}{\beta_+^2}$$
, beace  $0 = \frac{\varphi^{-1}}{\varphi^{-1/2}} = \frac{1}{2} \frac{1}{\gamma - \beta_+^2}$  and therefore

$$\alpha = \sum_{\gamma = \beta_s}^{-1} = \sum_{\gamma = \beta_s = \gamma}^{-1} \frac{\gamma - \beta_s}{|\gamma - \beta_s|} = \sum_{\gamma = \beta_s} (\gamma - \beta_s) \cdot b_s$$
 where  $b_s$  is positive

We need the geometrical representation of the complex numbers in the pane. In the sum is equal to 6, every component of this aum is equal to 6, every component of this aum is equal to 6. Let G be an arbitrary straight line passing throughly. The components of  $\gamma = I_1 - I_2$  orthogonal to G form a sum equal to zero, bears either the components are all equal to zero or there are components with different again in the let case the points  $B_1$  are all attented on G in the 2nd case, there are roots of  $\phi$  on both adea of G. In no case there are mote of  $\phi$  on no ado of G only Let now P be a concex polygon or idea, a the roots of  $\phi$ . If  $\gamma$  is outside if P we can draw a straight not that intersecing it through  $\gamma$ . Hence P and therefore all the roots of  $\phi$  are idealed on the same sub of G. Hence  $\gamma$  is not a root of  $\phi'$ .

Let  $\Gamma_0$  be the smallest claves proposited and the rots of  $\phi$  (The reader may prove that such a polygon exists and is image,  $P_1$  the direction of polygon defined by  $\phi'$ .  $P_2$  the smallest polygon containing the roots of  $\phi$ . The physons with higher naives are included in the processing  $\phi'$  degenerates to the point  $\frac{1}{2} = \frac{1}{2} \sum_{i=1}^{n} \frac{1}{2} \sum_{$ 

of a mad for the same reason to the centre of gravity of the roots of quantities of the roots of each derivate.

## § 5. INTERPOLATION

[5/1] Let

$$\beta_1, \dots, \beta_{n+1}$$
 (6, 1)

be n+1 i flerent e ements of an artitrary Seld K, and let

$$A_3 = \{A_{n+1} = \{5, 2\}$$

be n + 1 arbitrary elements of K



#### INTERPOLATION

We said to the foot a polynomial to of K folia that

$$f = v_s$$
 for 1 and degree  $\leq a$ 

I cold to a component satisfy

$$\sum_{i=1}^n B_i f_i = h_i$$

The between nect of the a set of a a bound of a set of a set of the part of the set of t

Let  $I_{\ell}$  be the solution of  $|\nabla_{\varphi}| = 0$  ,  $|\nabla_{\varphi}| = 1$ , then  $|f| = \sup_{i \in I} ||X_i f||_{L^2}$ 

is the secution for arbitrary techs ents. That to a second the

where year II to s. So we get by on year in ula country diten

$$f(x) = \sum_{j=1}^{n} \frac{\lambda_{j}}{1 - \beta_{j}} \frac{\lambda_{j}}{j}$$

$$(1)$$

By Lagrange a formula the process of interports a law seasons of [5/] in the most of my steeped general measure, but the formula we not convenient for practical calculation. It is essue to introduce the convenient of the presentation of the convenient of t

$$f(x) = \gamma_1 + \gamma_1 + \gamma_2 + \beta_1 + \gamma_2 + \beta_1 + \cdots + \beta_2 + \cdots + \gamma_2 + \cdots + \gamma_3 + \cdots$$

Here is  $\gamma_0 = \ell = \lambda_1 - \gamma_1 = \frac{\lambda_2 - \lambda_1}{\ell_2 + \lambda_1}$  and we may successively

enfoulate the configurate y. It is convenient to arrange the set united in the following manner.

Let file be defined by fire wire, and for kml n

$$f_{\theta}(x) = \frac{-t_{\theta-1}(x) - t_{\theta-1}(d_{\theta})}{x}$$
,

then  $f_1(e) = y_1 + y_{1-1}(e - \theta_{-1}) + \cdots + y_{n-1} = \cdots + y_{n-1}$ 

ALGEBRA

Hence  $f_{+}(\beta_{+-1} = \gamma_{+})$  We have therefore to calculate the values  $\{k, m\} = f_{+}(\beta_{-})$  for  $k = 0, .m, k < m \le n + 1$ 

by 
$$\{k, m\} = [\{k-1, m\} - \{k-1, k\}] - \beta_m - \beta_k$$

and {0, 0} = t. We call plate the values column vise in the following achieve

$$\{0, m\}$$
  $\{1, m\}$   $\{2, m\}$  ...  $\{n, m\}$   $\{k, 1\}$   $\lambda_1$   $\lambda_2 = \lambda_1$   $\lambda_3 = \frac{\lambda_2 - \lambda_1}{\beta_3 - \beta_1}$ 

$$\{h(-it)\}$$
  $X_3 = \frac{\lambda_3 - \lambda_3}{\lambda_3 - \beta_4} = \{1, 3\} - \{1, 2\}$ 

$$\{k,n+1\}, \lambda_{n-1} = \frac{\lambda_{n-1}}{\beta_{n+1}} = \frac{\lambda_1}{\beta_1} = \{1,n+1\} = \{1,2\}, \quad \{2-1,n+1\} = \{n-1,n\}, \\ \beta_{n+1} = \beta_1, \quad \beta_{n+1} = \beta_n, \quad \beta_{n+1} = \beta_n,$$

The first elements of the 1 fferent commune of this scheme form the set Yo Yir. You of the coefficients of 40. This scheme is cased for calculation than Lagrange a formula

[5/8] The reckening can further be sumplified if the elements  $\beta_1, \dots, \beta_{n+1}$  are equidistant, i.e. if  $\beta_{n+1} + \beta_n = \Delta_n$ 

for every k ; then

$$\Delta_{d} = \{k \mid m\} = \{\{k-1, m\} - \{k-1, k\}\} = \{m-k\}$$

$$= \frac{1}{m-k} = \sum_{k=1}^{m-1} |\Delta_{k+1, k+1}|,$$

where  $\Delta_{i,j} = \{k-1,j+1\}$   $\{k-1,j\}$  is the difference of two consecutive remember to the precedure estimate  $\Delta_i \{k,m\}$  is the mean table of the differences if consecutive elements of the rows m to k in the column k-1

We will now transform the scheme for these cases in such a manner that we have to calculate the differences of consecutive elements only and not the mean-values.



For this purpose we introduce the polisions usual in the calculus of differences.

Let 
$$\Delta_x$$
,  $u=x+\beta_1$ , then  $\Delta_x(u-k-1)=x-\beta_x$ 

Let 
$$F(u) = \int x_1 = \int \Delta_x u + i t_1 = y_0 + \gamma_1 u \Delta_x + \gamma_2 u - (u - 1) \Delta_x^2 + \dots + y_n u$$
,  $(u - 1) \dots (u - n - 1) \Delta_x^2$ 

and let  $\Delta t(x_i = f|x + \Delta_{xi} - f|x_i = 1, u + 1, -1, u$  then

$$\Delta f(x) = \Delta_{+} \left[ \gamma_{1} + 2\gamma_{2} u \Delta_{+} + \cdots + n_{1} u - u - 1 - u - u_{+} \cdot 2i - \Delta_{+}^{-n-1} \right]$$
  
 $\epsilon_{12} = (u+1)u(u+1) \cdot (u-k+1) - u(u-1) \cdot ... (u-k) = (k+1)u(u-1) - (u-k+1).$ 

Let 
$$\Delta(\Delta f | x) = \Delta^{\otimes} f(x)$$
,  $\Delta(\Delta^{\circ} f | x) = \Delta^{\circ \circ \circ} f(x)$ ,

then we get by repetit on of the proced tre

$$\Delta^{\frac{n}{2}}f(z) = \Delta^{\frac{n}{2}}\left[2\gamma_{\frac{n}{2}} + 2^{-1}\gamma_{\frac{n}{2}}u\Delta_{\frac{n}{2}} + \cdots + n^{-n} - 1^{-1}\gamma_{\frac{n}{2}}u^{n}u + 1\right] - u - n - 3^{-1}\Delta_{\frac{n}{2}}u^{n-\frac{n}{2}}$$

$$\Delta^n f(x) = \Delta_n^n \cdot n!$$
 yn

For abbreviation we shall write  $\Delta^2 f/I_0 = \Delta_0^2$ . Then

and

$$f(\theta_1) = \gamma_0$$

$$\Delta | = \Delta_{a71}$$

$$\Delta \{-\Delta_s k | \gamma_s$$

$$\Delta_i^n = \Delta_i^n n \mid \gamma_n$$
 (for  $i = 1, 2,...$ ) holds.

So we get Newton's formula

$$f = f(n_1 + \Delta \{ u + \frac{1}{2} \Delta \} u + 1 + \dots + \frac{1}{n} \Delta_n^2 + u + 1)$$

$$= f(s_1) + \frac{\Delta_1^2}{\Delta_2^2} \left( s^{-1/2} + \frac{1}{2} \right) \frac{\Delta_2^2}{\Delta_2^2} \left( s + s_1 - s + s_2 + \cdots + \frac{1}{n} \right) \frac{\Delta_2^2}{\Delta_2^2} \left( s + s_2 + \cdots + s_n + \cdots + \frac{1}{n} \right) \frac{\Delta_2^2}{\Delta_2^2} \left( s + s_2 + \cdots + s_n + \cdots + s$$

$$x = \mu_{-1}$$



The elements A; can be calculated very easily by the following scheme —

The degrees of f(x),  $\Delta f(x) = \Delta^{n}f(x)$  are decreasing and the last one is a constant so we can use the above achieve a no for extrapolatica to get the value of f(x) for every arbitrary one geal value of f(x) that means for every since  $f(x) = d_1 + \kappa \Delta$ , where k is an arbitrary integral number

From  $p^{\ell}$  Let f(x) be of legree 1 and let  $\ell$  be 1  $\ell$  2 = 4  $\ell$  0 = 5,  $\ell$  4) = 1  $\ell$  5 = 2. In order 1 get  $\ell$  1 we use the exhemo, order using at first the numbers ab so the ditted line from the right to the right and then the numbers below the datted line from the right to the left.

Hence f(0) = 33.

## PART V

MATRICES. RESULTANTS.

## § 1. Marnioles

In the first part of these lectures matrices have been used to solve [1 1] evalues of linear equations on him 14 and 15 a few properties of matrices have been discussed. We will now consider the matrices in a more systematic manner.

Let K be an arbitrary field (see Part 11 2 let 0 be the notletement 1 he the unitelement of K, and let

be arbitrary elements of K. A whome of m. raws and a columns "

$$\begin{pmatrix} (1 - a \cdot 1) \\ a_{1-1}^{\alpha} & a_{2}^{\alpha} \end{pmatrix} \qquad (-a_{2}^{\alpha}) = \lambda \qquad (1.2)$$

has been called (see Part I 6) a matrix and  $\tau$ ; its elements. If  $n \to 0$  this number will be suid to be the direct of A, the general case can be reduced to the case m = n

Let A, B C, be matrices of degree C and let the elements be denoted by the corresponding small latin type as an (1 a. The galf w of matrices is given by

$$A+B=P$$
, where (1, 3)

a + b = f for i = 1 n k = 1 n

The commutative law 
$$A + B = B + A$$
 1.4)

and the associative law 
$$(X+B+C-A+(B+C))$$
 1.5) hold for this addition of matrices.

hold for this addition of macricos.

Let r be an arbitrary element of the held h, then we define the product  $e A = f(c \circ d) \cdot 1 \qquad (1.6)$ 

e . we multiply every element of A with and get a new matrix a A of dogene s. Then the

<sup>\*</sup> Instead of brackets sometimes vertical double tare are put

72

ALCHURA

hold, For c=0, we get

the attempt of Oberg quality and

If we define the subtraction by

$$A = B = A + (-1) B, \qquad (1, 10)$$

the west to be an analytic part of an analytic part of the many of the state of the

When Matter is a known furth 1. Much be proved to be created by a finite of the first transfer of E  $\alpha$  for F (r,s) is defined by  $\alpha$ ,  $\alpha$ ,  $\alpha$ ,  $\alpha$ ,  $\alpha$ .

$$r_k^*(r, s) = 0$$
 for  $(r, k) = (r, s)$  (1, 11)

thon

$$\underset{t \downarrow k}{\mathbf{X}} \in \mathbb{F}(t, k) := A \tag{1, 12}$$

the send the part found only forery ; a equal to 0. Here the s<sup>2</sup> restricts were if a lept 10-5 form a few of M, and M can be 0. set to a series, if each s<sup>2</sup>, see Part 1. 4. Part II (> 1). By the accessory to 1 for the industribution to the sum of the part of the industribution to an arbitrary biomorphism.

[1 2] To Part I II the most points as of material of M has been defined by

$$A B = D, d = \sum_{j} a_{j}^{*} b_{j}^{*}$$
 (1, 18)

and he need the new ABC - A the

(3, 14)

have been president lead (1.1 the

dark too have S+LC AC+BC

$$C(A+B) = CA+CB \qquad (1, 16)$$

follow directly.

Minister for a 17 we Post II | 1 ? The mater may prove as an error so that the ring is non-ministrate except when no 1. The ring has as the unitelement, the matrix

$$F = a_k^* \quad \text{where } a_k^* = 1 \text{ and } a_k^* = 0 \text{ for spik} \qquad \qquad 1.10$$

MATRICUA

73

Then

BARAKSA.

(1, 17)

and

 $0 \pm 0.0 \pm 0.0$ 

(1, 19)

for every matrix A and B and A and A are A and B and A are A are A and A are A are A and A are A and A are A and A are A are A and A are A are A are A and A are A and A are A are A are A and A are A are A and A are A are A and A are A are A are A are A are A and A are A are A are A and A are A are A are A and A are A and A are A

To every restrict to the every many to the many distance Part I § 16

(1, 19)

On the ther hand, to are more to the expectation of a matrix to the consults or equate the trace of the consults or equate the trace of the matrix of the consults of the cons

Residence to product the second of the second of Known to the second of the second of

$$det(eA) = e^{-t} det A$$

En mout the wife every mark H.

AN TO YARD PARK TO A CONTROL OF THE OWNER OF THE PARK TO A CONTROL OF THE OWNER O

$$L_{ij}^{*} = abj^{*}$$
 (1, 22)

he the cofector of all two Part I p. 421,

• Physical Response for the second se

then  $\Sigma a(b) = 0 = \Sigma a(b)$  for  $i \neq k$ 

$$\Sigma = \ell_1 = 1 = 2 + 1$$
 hills see Pirt I i and (  $\ell$ )] Here

$$AB = E - BA \tag{1, 20}$$

hada. The matrix B the elements of which have been foliated by 1 22) is said to be the once or of A aut will be 3 in declary.

$$8e^{-}AA^{-1} = E = A^{-1}A$$
 (1, 23)

There exist therefore an inverse mater of and only if the determinant is different from 0

Let dat A = 0, then from

$$AX = B$$
,  $YA = B$  it follows that  $X = A^{-1}B$ ,  $Y = BA^{-1}$ 

Hence the markers with non-versity it rought from each in which the two overset at and give a unique resent. The area on the result of the arm them there will be non-versity determined at any to equal to 0. Of coursely of the arm of the arm of the arm of the arm of the reserve of the arm of the a

[1/8] Where to consider a with a statement to every matrix to the extra local terms from the process of the avertors were known from the form of the extra local transfer and transfer and the extra local transfer and transfe

$$a(x_1 + a)x_2 + ... + a)x_n = a'_1$$
 (1, 24)  
 $c = 1, ..., n$ 

By this transformation the unit vectors over l'art lip " of her production formed to the event of the life of the end of the life of the l

Let 
$$\begin{pmatrix} x_1 & 0 & \cdots & 0 \\ x_2 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \ddots \end{pmatrix} = \cdots \tag{1, 25}$$



then (1,24) becomes

$$A(s) = (st)$$
 (1, 24)  
 $(s) = A^{-1}$ 

the forms is all fill and a property the transfermation by which the basis of the a vectors, and because it is the unity of many because transformed to the vectors. If offer

$$1 \leftarrow k \quad r \Rightarrow 4k \quad q \qquad r = 18 \quad \text{and } d \in 13 \Rightarrow 1, \qquad 1 \Rightarrow 20,$$

then

$$AB(y) = B(y)$$

$$B^{-1}AB(y) = (y'),$$
 (1, 27)

By (1 26) a mean of the respective of the very fithe correspond the vectors

$$(B_1) = (b_1, ..., b_n^n)$$
 (1, 29)

I the engagement the exchange from a time and environ y each basis of he vect requestions of a matrix B with a name thing determinant. The vectors of a matrix B with a name thing determinant. The vectors of the By become transferred by 1 27 on the same manner on the anit vectors become transferred by 1 24 to \$\frac{1}{2} \frac{1}{2} \frac{1

is said to be the fr million of A by B . As

hold the transfering of a first restrict 1 from a large see fact II I I and from (170) to fill we that the maters of the constraint the transfering of a generated by inflerent materies of the name of the sectors passe. If we require the university by the basis formed by the vectors of the vectors of the vectors of the vectors of the energy of the sectors of the energy of the energy of the energy of the energy of the sectors of the energy of the sectors of the energy of the

Some purbounderate it a make that he transform of a by b. !

41.4

## 12. TRANSFORMATION OF A INTO A NORMAL-PORM

[2/1] It A is a matrix call of a for K at V be a sufficient at a form a command V. I the set weeder I a so the season of the set of the season of the season

The coud to no

.1

$$a \uparrow b \uparrow + a \uparrow b \uparrow + ... + e^+ b \uparrow = A, b \uparrow$$
 (2. 1)  
for  $k = 1$  (fixed

form a system of a basis, make the conquestions with the matrix

$$\mathbf{A} = \mathbf{\lambda}_1 \cdot \mathbf{E} \tag{2.2}$$

Therefore are a star at the color, a familiarly f

$$\det \left( \mathbf{A} + \mathbf{\lambda}_{+} \mathbf{E} \right) = 0,$$

is all promise to a b , we field to now as mappoint to onto a

the n roots  $\lambda_1, \dots, \lambda_n$  of (2,3) and

$$\chi_{\pm}(x) = \Pi (\lambda_{\pm} - x),$$
 (2.4)

[2/2] for a case dead invest, it is find an arritrary single maters and to projecte it a maters twice and the last transfer the multiplication in the usual manner.

$$A = E A^{\dagger} = A, A \cdot A A , A^{(+)} = A^{\dagger}A$$
 4.5)



and if dot \$ 500 \$ 7 m \$ \$ , take \$1 we that

$$\lambda^{*} \lambda^{*} = \lambda^{**} = \lambda^{*} \lambda^{*}, \qquad (2, 6)$$

The powers of A are commutative matrices.

Let  $\phi(A)$  the matrix

$$\phi(X) = \sum_{i=1}^{n} \alpha_i(X_i) + \alpha$$

and  $\phi$  will be said to a montegen on the fit of the very late  $\phi$  be a polynomial to V ,  $v = -\phi + 1$  and  $v = -\phi + 1$  follows that

$$\phi(\Lambda) \psi(\Lambda) = \phi(\Lambda) = \psi(\Lambda) \phi(\Lambda), \tag{2.6}$$

There the integral function for the matery to come a formal of the community of the communi

therefore cann the independent, and there must be a polen new to condense of the degree of not with the property that

$$\omega(\Lambda) = 0$$
.

vectors my -1 . I which are leponiers but every amin's moves

of them is independent. Then an equation

$$a_1(\beta_1) + \dots + a_r(\beta_r) = (0)$$
 (2, 0)

bod of response for the after (21) by A

Wo get

$$u_1\lambda_1(\beta_1) + \dots + u_r\lambda_r(\beta_r) = (0)$$

He part of + + + ,

The control of the second of t

to the meanth to the property of the second different A case in

transformed into the diagonal-matrix

$$\begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \dots & 0 & \lambda_n \end{pmatrix}$$
. (2,10)

From Record to the first the perdent, home because Record to the first the contract of the state of the first the section of the first the section of the se

If m rets are not u. I f. mt to con larv down not hold in cvery

esco. Eq. Let 
$$A = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$$
 then  $\chi_{A} = 0 = 1 + 2^{A}$ ,  $\chi_{A} = \lambda_{B} = 1$ 

V:  $\approx \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}$  has the rank 1, when 2.1 has into me solution  $(x_1, x_0) = (1, 0)$ .

Σ<sub>0.1</sub> (B<sup>-1</sup> A B)<sup>1</sup> = B<sup>-1</sup> Σ<sub>0.1</sub> A'B bolds-

Theorem Transform a control in to by how an equal to A semilar theorem in the factor of the exercise years.

There are I transfer not in the area or positive same characteristic polynomial

Proof As 
$$bet^{A}B$$
 for  $B$  ,  $a = 1 + B + B$  ,  $a = 1$  dot  $(B^{-1}A + B)$   $a = dot A$  holds  $B^{-1}(A + aE)B$   $a = B^{-1}A + B + aE$  ; hence dot  $A = aE$  is  $B^{-1}(A + aE)B$   $a = B^{-1}A + B + aE$  ; hence  $A = aE$  is  $A = aE$ .

The two last theorems have a certain right in his enthron. It will be proved after in the excess mark your of the histories of the histories. Now we will consider only that can while the training polynomial are all different.

As the matrices 
$$A = \{ \{1, \dots, 1, \dots, (A-\lambda_1 E)(\beta_1) = 0 \} \text{ As the most rices}$$
 A type are contrated to  $\lambda_1 = 0$  . If  $\lambda_1 \in F$  ,  $\lambda_2 = 0$  .

The vector franchista transfer by the transfer of rank to

Rence (see Part 1 5 14)

$$\lambda_{A}(A) = 0,$$
 (2.11)

In the formula the ster ks tenne certain of marks which test be consulated parameters and A a a new rest of discount of

 $\Lambda_{\Sigma}\chi_{A}=x^{*}\times\chi_{B_{2}}\chi_{A}$ ,  $x\to \infty$  is  $\chi_{A}=x^{*}$  in the transform  $\Lambda'$  into

$$B'A'B'^{-1} = \begin{bmatrix} \lambda_0 & \dots & \gamma \\ 0 & A'' \end{bmatrix}$$

where A" is of degree a -2, and if

$$B_{+} = \frac{1}{1} = \frac{1}{1$$

The first rive of the state of the state of the countries of the proposition of the proposition was get.

A "nearmant fagure r  $\chi_{A}$  results on which is not described by  $B_{a}AB_{a}^{-1}$ . Hence

$$\begin{aligned} \mathbf{A}_{-1} &= \mathbf{A} \\ \mathbf{A}(\beta_{\pm}) &= \lambda(\beta_{\pm}) + c(\beta_{\pm}) \\ \mathbf{A}(\beta_{\pm}) &= \lambda(\beta_{\pm}) + d_{\pm}(\beta_{\pm}) + d_{\pm}(\beta_{\pm}) \end{aligned}$$

$$A(B_r) = \lambda(B_r) + h_1(B_1) + ... + k_{r-1}(B_{r-1}),$$



Therefore

$$tN = t\ln - \lambda_1 t = 0$$
 (2.15)

$$(A = AE)(B_0) \Rightarrow c(B_0), (A = AE)^n \cdot (B_0) = c(0)$$

$$(X \cap A) : C = \{I_1C = A \mid I_1 \cap I_2^{-1} \mid X \cap A\} \quad \forall i, i = \{I_2C = \{A \mid X_i^{2} \mid I_{i+1}\} = I\}$$

$$A = AT + B = 00$$
,

Homes for every voctor of all the vectorspace V greented by of the

$$(A - \lambda E)^* (a) = (0)$$
 (2, 16)

to be The vestors of a restrict with non-variating lo-

We want to prove new that every restor to, for which

$$(A - \lambda E)^{\pm} (\gamma) = (0)$$
 (9, 17)

In the tree range to V. If there we used because his a vector of V. I. c.

$$A(\gamma) = \lambda(\gamma) + h_1(\beta_1) + ... + h_r(\beta_r).$$

An Int of the pare appeared to be independent there is also minimised, if which there exists form the z + 1/1 int rown

Thou

that therefore  $\chi_{A}(x) = \lambda - x$ ,  $\lambda_{A}(x) = \lambda_{A}(x)$ . The vectorspace 2, 17) is therefore composed of all vectors which satisfy (2.1) for my exponent f. As A = AE) in  $f = a^{2}$  satisfies (2,17) if a considerations we get the following theorem.

Theorem If A is a root of  $\chi_{A}$  , a) of order r them the vector seast stying (2, 17, for any q form a vectorspace V of rank r. Each vector i) of V satisfies 2, the and is teamsterized by A to a vector of V. If the first r on amount of B. form a spinor y bosen busin if V, then (2, 14) heads

[2/5] Leb

$$\chi_{_{A}}^{-}(z) \rightarrow (\lambda_{1}-z)^{r_{1}} \dots (\lambda_{n}-z)^{r_{m}},$$
 (2, 18)

where \(\lambda\_1\), \(\lambda\_m\) are liftered. Thereby \(\lambda\_i\) there corresponds a vector space \(\lambda\_i\) of rank \(\lambda\_i\), so that for every vector \(\sigma\_i\) of \(\lambda\_i\), the

equation  $(A-\lambda_*E)^{r_*}(a_*) = (0)$  holds and  $(a_*)$  is transformed by A into a vector of  $V_*$ . To prove that the vectorspaces  $V_*$  are independent, we have to use the following lemma

I wrome. This has the profession and the professio

From The lemma a true if m=1 and m=2 see Part II 4/5 i, let it true for m=1 and therefore the h of of  $\phi_1$ ,  $\phi_{i+1}$  is  $a_i=a_{i+1}$  for  $a_{i+1}+\cdots+a_{i+1}$  for  $a_{i+1}=a_{i+1}=a_{i+1}$  by the h of  $\phi$  of  $\phi_1$ ,  $\phi_{i+1}$  is the h of  $\phi$  and  $\phi_{i+1}$  there are  $a_i=a_{i+1}$  and  $a_{i+1}$  is  $a_{i+1}$  and  $a_{i+1}$  and  $a_{i+1}$  is  $a_{i+1}$  and  $a_{i+1}$  and  $a_{i+1}$  is  $a_{i+1}$  and  $a_{i+1}$  and  $a_{i+1}$  is  $a_{i+1}$  and  $a_{i+1}$  and  $a_{i+1}$  and  $a_{i+1}$  is  $a_{i+1}$  and  $a_{i+1}$  and

The press I ng learns with walp not to the polymon as

$$\phi_A(x) := \chi_A(x) : (x - \lambda_A)^{T_A}$$

As these polynomials are restively prime, there exist m polynomials  $\eta_1$  (a) satisfying

$$q_1(x) + ... + q_n(x) = 1$$
 (2, 19)

 $\eta_{+}(x)$  a divisible by  $\phi_{+}(x)$   $i=1,\ldots,m$  and therefore

by 
$$(a-\lambda_0)^{\ell_0}$$
 for  $k \neq i$ .

Hence

 $n_0 \in A$   $\alpha_0 = 00$  and from 2/14 it

2, 24

follows that

 $\psi_i$  (A)  $\langle \alpha_i \rangle = \langle \alpha_i \rangle$ ,

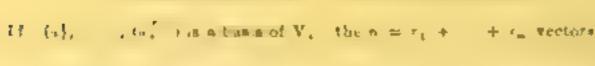
If a section we I the fift for sect reperce V, untially

$$(a_1) + ... + (a_n) = 0$$

we get by mustiplestion with the matrices of (A for every t

$$\Psi_{1}(A)(a_{4}) = (a_{4}) = \{0\}$$

Hence the vectorspaces V. . . . V. are nispendent.



$$\{a_i\}, \dots, \{a_i^{f_i}\}, \{a_i\}, \dots, \{a_n^{f_i}\}$$
 (2, 21)

for n a basis of the space of the n victime. The matrix C whose relumns are the vectors (2.21 has therefore a determ number 0. The vectorspaces V are invariant for the transformation 4. Hence the neverterspaces generated by

$$\{e_{i_1+i_2}, \dots, e_{i_{n-1}+i_{n-1}}\}$$

are invariant for C A C . The matrix has therefore the form

A.

where A, as a matrix of degree r, and its characterists planers is equal to (t,-x). I see some stanted outside he square is equal to 0. As it has been proved in [2.4], the equation (x-k,1) = 0 had for every vertex a, lof V. As the vectors (2.21 form a basis of the total vectors every vector vector a is the sectors (x-k,1) = 0.

$$\chi_{A}(A) = 0,$$
 Hence  $\chi_{A}(A) = 0$ .

The east are of these count local mans given by the fellowing theoreta.

There we to the character of and r. . = 1 to the recommon of C re bus a

the form (2.2) A sale of this which are the polynomial

[2/0] If the bosos of V is repaid by another have the matrix A, f (2/22) we be transform as and up to the matrix of A, \$\phi\$, remaining and tered for put the matrix 2 of the rime form at an therefore on the necessary to transfer the matrix \$\psi\$, and there is the form of A is equal to

$$\chi_{\lambda_{-}}(x)=(x-\lambda)^{+}.$$

forth to be out to the the scan speak of the which

$$(\Lambda + \lambda E_1^{-1}(\beta) = (0)$$
 (2, 23)

It is the state of the state of

$$q = \exp(\beta),$$
 (2.24)

Lot 0 < p < q , e | 0 , then

 $\exp\left(c\beta\right) = q$ ,  $\exp\left[\left(A - \lambda R\right)^{p}\left(\beta\right)\right] = q - p$ , and

$$|A + A| = r + |A| \qquad \text{on } A_1 + A_2 \leq \exp(A_1)$$

From the country of the vectors for which cap  $B_1 \le C$  to the contract of the vector spaces

$$W_1, W_2, ..., W_r = V$$
 (2, 24)

evers center passe in reliabled in the subsequent space at many also be identical to it. For every a < r the vectors

form a victoria a contribution with in W. The most of the torigine activity will be leaded with the Besselle to the segment of the victoriances.

$$\mathbf{W}_{\mathbf{r}_1,\mathbf{r}_2+\mathbf{p}_1} = \mathbf{W}_{\mathbf{r}_1,\mathbf{r}_2+\mathbf{p}_1,\ldots,\mathbf{r}_n} = \mathbf{W}_{\mathbf{r}_1,\mathbf{p}_1} = \mathbf{W}_{\mathbf{r}_1}$$



every vectorspace is in lated in the subsequent spaces. Lit

be the rank of the vectorspaces

$$W_{1+r-1}, W_{1+r-1}, W_{1-d-1}, W_1$$

then there exists a basis

$$(\beta_1^1)$$
, ...,  $(\beta_{n-1}^1)$ ,  $(\beta_{n-1}^1)$ ,  $(\beta_{n-1}^1)$ ,  $(\beta_{n-1}^1)$ , ...,  $(\beta_{n-1}^1)$ 

of W1 with the property that

$$(\beta_1), \dots, (\beta_{k_n})$$

$$(A - \lambda E)^{\perp}(\beta)^{\perp 1} = (\beta!)$$
 3.40)

and we define ( $\beta$ )), for  $1 < k \le s + 1$  by

$$(A + \lambda E)^{a+a+b} (\beta I_{i+1}^{a+b}) = (\beta I_{i}^{a})$$
 (3, 26)

From 2 26 tifst synthat 2 26 by defect = 14

Let the weet less of the arranged in a triangular actionic

$$(B_{+}^{\dagger}) = (B_{+}^{\dagger})_{1} = (B_{+}^{\dagger})_{2} = (B_{+}^{\dagger})_{2$$

$$(\beta_1^q)$$
 , ...,  $(\beta_{r+1}^q)$  ,  $(\beta_{r+1}^q)$  , ...,  $(\beta_{r+1}^q)$ 

$$(\beta_1^*)_+ \dots, (\beta_{r-1}^*)$$

If we multiply a vector of this a home from the off one with the matrix (A-AE we get the vector just up yout

$$(A-\lambda E)(\beta_{+}^{n+1}) = (\beta_{+}^{n})$$
 hence

$$\Lambda(\beta_+^{n-1}) = \lambda(\beta_+^{n-1}) + (\beta_+^n)$$
 holds.

The volt reported by the sectors of an arbitrary form may the C , a therefore my mount under the transformation by A. These will be

are independent, ris., from

$$(0) = \sigma_1(\beta_i^1) + \dots + \sigma_k(\beta_k^k)$$
 it follows that 
$$(0) = (A - AE^{-k+1} - A^k) - \epsilon G^k$$
 bence  $C_k = 0$ .

The first of the first bearing of the basis of the basis.

where the diagonal chimer is are 1, in the adjacent parallel and the chaments are 1 and the others in the are 0. The degree of 2 28) is equal to a + 1 where rise given by 2.2. In get the a remaindering of the tempelormation A we have to prove that the vectors of 2.27, are independent and form a basis of 3. The vectors of the first rise are independent and form a basis of 3. The vectors of the first rise are independent and form a basis of William vectors.

ingle prestors with  $\lambda = \lambda E$  that  $\sum_{i=1}^{n-1} -\beta^{n} = -1$  but as the vectors  $(\beta \})$  are

independent, 
$$0 = c_1$$
 , holds, and therefore  $\beta^{ij} = \beta$ 

Hope the set is 4 the two tra fore are independent. By mathomat along so on 1 to the " that the vectors 2 27 are adependent

$$\chi = (1 - \beta)^p = \sum_{i=1}^{n_p} I_{i-1}^{-1}$$
), as it is a vector of  $\mathbf{W}_{1+1}$ 

Let 
$$(B^{\frac{1}{2}} - \frac{N}{2})^{\frac{1}{2}} = B - \epsilon - \rho$$
 then  $A \rightarrow Ve^{-}(\gamma) = 0$  bears  $(\gamma)$  in a

vector of Walls

$$\langle i_i \rangle = \sum_{i=1}^{k_0} k_i \cdot \beta_i \rangle$$

<sup>\*</sup> The safe repair revers at these combinements a section to on on accordance



 $\beta^2$  in therefore dependent on the vectors of the two first rows. Hence these vectors from a basis of  $W_2$ . If  $f : 2^2 = 4$  vectors  $f : 3^2 = 1$ , then  $A^2 : 2^2 = 1$ . If  $A : 3^2 = 1$  belongs to  $A^2 : 3^2 = 1$ . Then  $A : 3^2 : 3^2 = 1$ .

The vectors  $(\beta_1), \dots, (\beta_{n_1}, \beta_1), \dots, \beta_{n_2}^*$  from therefore a beam of the vectorspace  $W_{n,1}$ .

By mathematical major in " at forms that there exert the heat k rows form a base of  $W_{k+1}$  so we get that for k=r, the vertex 12.27 form a base of the total vectorspace V

We arrange this waste a ring to the cocations, and we considere in such a manner that the vectors of the base books on the lore. Then A will be transformed into

The transformation of the subspaces U is given by A. Isomer this constrains have the a small real 2.34). The legacing the mark A campant to the length of the corresponding regions of 2,27). In it discusses the a market of the vote part Queen in the member form 2,271. 3.24, the ranks of the vote part W, W, one easily as found out the ranks of W, being the notices of the matrices A', of six a the digres is not sess than A.

In the general case we wan transfer belong to I Treate one of nature to In the general case we wan transfer on each matrix A. If 2 22 no he normal form. By these considerations the new way the rem has each proved.

· See y. St. Jostophi.

maters A. Then A can be transferred into an a rational (2, 22), the material A<sub>1</sub>, ..., A<sub>n</sub> have to form, which and the form A<sub>1</sub>, ..., A<sub>n</sub> have to form, which of the n backle form 2 is not be a solution of the n backle form 2 is not be a solution in the transfer arbitrary.

# 5.1 SOME CORRESPONDED FOR COMMENSOR OF THE CHARLE PROJECT POLICE MALE.

the comments was been restricted and he change of more easily the light the case to a more employee range.

[4,1] I to recent the process of the continuous state as a superficient to the process of the p

The regard for a Character of Character of South

$$\begin{pmatrix} \gamma_1, & \gamma_0, & \dots & \gamma_r \\ & \gamma_1, & \dots, & \gamma_{r-1} \\ & & \dots \\ & & \gamma_1 \end{pmatrix}$$
 (9, 1)

where  $r_1 \neq 0$  and the commute cover the diagonal are equal to 0. Here if 222 is commutered to 3 and there is no other matrix countries of 22, then the matrices to 1 as every other transforms that will not give a second the volume of the volume of the matrix.

If  $u_i = c$  on b of the matrices V is I degree 1. Hence  $0 = u_{-1} = -1$ , in the assets went replace it lemical with V. There are no spaces  $V_{1+1}$ , every basis of the vectors are leads to the prime form. The normal form is therefore commutative to every matrix with I term name  $\frac{1}{2} 0$  of course the adaption matrix with the diagram, elements  $\frac{1}{2} V$  in the general case the investigation of the same strength random this of the case v = c the investigation of the matrices commutative to be pursue form as more complicated and we



will not go into the details here. The more general problem of finding out the material committative to an inhibitive and has been at A. committee; be reduced to that problem in the of A into Direct commutative of A in and C B C<sup>-1</sup> are commutative too.

Thomselvis to a roo of the barieter to p yours as \(\chi\_1 \text{tx}\), but [1/2] it may be that A a a cut of a polynomial for ever police ogree. Let

$$\chi_A (a) = (A - \lambda_A E)^{-1} (A - \lambda_B E)^{-2} \dots (A - \lambda_m)^{-m}$$

To  $V = V_1 F^{-1}$  there excess wis the mass materials  $V_1 = V_2 F^{-1}$  the normal from 2, 22, and the mass over that the from 2.28 is appeared by diagonal matrices  $V_2 = V_2 F^{-1} = V_3 = V_4 = V_4$  be the degree of  $V_2 = V_3 = V_4 = V_4$ 

$$(\Lambda - \lambda_1 \mathbf{E})^{-1} (a_1) = (0)$$

Cherentending to the detailed on which the later of a least then

$$1 \le r_* \le r_* \qquad \text{and}$$
 
$$+ (A + \lambda_* P)^{r_*} (s_*) = 00 \qquad \text{hold},$$

for every a tire a let be see a particle problem to A.

Let 
$$\phi(x) = (x - \lambda_1)^{-1} \dots (x - \lambda_n)^{-1}$$
, (3, 9)

then

$$\psi(A)\cdot(B)=(0)$$

hat to for every vector and of the section and marefus body for every wester of the vectorspace.

Hence 
$$\psi(A) = 0$$
,  $(0, 5)$ 

To prove that A is not a root of a polynomial who has not divided by  $\psi_{+}r_{+}$  we consider the case  $\chi_{+}r = \lambda - \epsilon F$ . Let  $\phi_{-} r = x - \lambda_{+}^{+} r_{-}$  between the legence of  $\Lambda_{-}^{+}$ . Then  $\phi_{-}^{+} x = \lambda_{-} + \lambda_{-}^{-} r_{-}$  where it is  $\phi_{-}^{+} x = \lambda_{-} + \lambda_{-}^{+} r_{-}$  tensor the transformation of the vectors is even in it is  $\phi_{-}^{-} x = \kappa_{-} + \kappa_{-} + \kappa_{-}^{+} r_{-}^{-} + \kappa_{-}^{+} r_{-$ 

and there exists Introduce a vector of an District P(y) = x > y. Hence  $\phi(A + i \neq 0)$  where force  $\phi(A \neq 0)$ 

follows directly that 
$$\nabla^2 = \begin{bmatrix} \mathbf{B}^{\#} \\ \vdots \end{bmatrix}$$
 and

Hen  $\sigma \omega(X) = \Omega$  if and on yet  $v(B) = \omega + v = \Omega$ 

The corresponding rule to it if A is companied in the two diagonal matrices. If therefore A has the form 2, 20,  $\omega$ ,  $\lambda$ ) = 0 hadrons 1 and only if  $\omega(A_1) = 0$  for i = 1, ..., m.

Hene we must be largible yearry and the first these common tems the first swing theorems has been proved.

Theorem  $at \lambda j = 0$  if and in y if e) is I was the by the polynomial  $at x_j$  which has been defined by A/2

[3/3] By the water and proposed of A the transformation generated by A becomes inequals defined only I the rows of the polynomial are all different. If there are equal roots, there are different to rough forms up I therefore if flerent non-mone spine transformations corresponding to the same characteristic polynomial.

To get the characteristic polynomial, it is not necessary to put the matrix in the mornal-form. As

$$\chi_{\lambda}(z) = \det (A - z E),$$

the coefficient of  $\pi^{+}$  in  $\chi_{A}(s)$  is

$$(-1)^{A} \sum_{n=1,1,1} A_{n=1,1,1}$$
 (8, 4)

where  $\Lambda_{n-k+1}$  denote the minors of A with n-k is we which are symmetric to the diagonal of A. The same are a direct to the transfermation;



the count epor ant fittees associate every set to k=n-1

We want a upper the theory to the business exhibition of a complex variable.

$$w = \frac{a + \beta}{a + \beta} \qquad w \text{ be recovery } \phi = \phi$$

We already to the energy of the first  $t = t_1 + tt_2 + t = t_3 + tt_4$ 

$$w_1 = u x_1 + \beta x_2$$

$$w_2 = \gamma x_1 + \Lambda x_2$$

As a commercial on the track of the section of the commercial of the section of

Hence  $\lambda_1\lambda_2=1, \lambda_1=re^{i\frac{\pi}{r}}, \lambda_2=r^{-1}e^{-i\frac{\pi}{r}}$ .

(II) 
$$\lambda_1$$
 )  $\lambda_2$ , it is remarkable from  $\begin{pmatrix} e^{i\frac{i\theta}{2}} & 0 \\ 1 & i \neq i \end{pmatrix}$  (if 5)

As a factor + 1 and a permutation of A by remain arbitrary, we can choose  $1 \le r_i 0 \le p < \pi$ .

2)  $\lambda_1 = \lambda_2 = \pm 1$ ,  $\kappa = \pm 2$ , by a suitable charge of the common factor =  $\pm 1$  we can arrange that  $\kappa = 2$ , but  $\sigma = \lambda_2 = \lambda_2 = 1$ 

There are two normal forms

$$\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix} \tag{3-6}$$

theory of functions. The classes of transformations with the norms form (3 5 are and to be first at especially for r=1 they are called rimple, for r>1 = 0 they are called the first nature (3, 6) denotes the identity the accommentation denotes a grand-nephron of the complex sphere being the only fixpoint and the other transformations of this class are called parabolic transformations, the only fixpoint being a finite point

## THEORY OF BLEMSNIAMS DIVISORS.

### Last Sharps or the ft had swing right a [ 6 2 ]

- a a very women. I to I be their perceptuals and into grad remotive (1 - 1)
- 2 If 5 is a fautor of a, ther

we are the again to have a familiary at the bore assumption

If you are trary every to be and a commet designification of then there on at in Normanata food in anticlying

$$f_{i_{1}} = f_{1}, N(a_{1}) < N(a_{2})$$
 (6. 8)

buse 1 th and have to a considered in art 11 ,4 4 and [1 i] of there out too the on those raps are

- the regular of the ottowed mainbour if we ofthe View of trey " " tontinefth regar - lant
- 2 The ring of the integer of orpox name value by face Part II t 11), where a and b are not ger. We do has North to by North and a his The and each there ag are +1 -1 +1 -1
- a the ring hard of the parts as a non-activerey indefinite a lover an array held & We network to fix every physicisms f(x) which is d flurent from the posyn more and if a madegree ( x) +1. The unities of this ring are the elements ‡0 of K

It has been proved in Part II of The factorization in S is unique and the had not be semented or I I in be easy read by

$$(a, b) = c + d b$$
 (4, 4)

where a, and d are elements of 5

As we see from the errors of [25] the hir find nich ments of 8 can be expressed by

$$(a_1,...,a_n) = a_1a_1 + a_2a_2 + ... + a_na_n.$$
 (4, 4)

Free see Green and red Storwin h (4 to (4, 2) and 4 3) leide prove that it a posts in it res age it function by a function beautistying the same conditions as 's so that 's is -1 if and only if it is a unity of S.



We now short in the method of except to a fact 1 to n [473] to altern with a country from 8 of a method and a now medical normal to be real normal and a series of the country as a series of the series we have to deal with the except to take from a real, the from a hard therefore to be reported by the age of method to be a top of the from a first of the posterior of the from a start of the from a series of the from the first different and a series of the from the series of the from the report of the from all the first different and from the first the first order of the from the right and a series of the series to made use oth interpretations at any step. The setting a series on have to made use the name as in Part I 1 15.

Diagona, in express 
$$D = 0$$
, where  $\omega = 0$ , so  $D = 0$ , for  $f \neq k$ ,  $df = 0$ ; (4, 5)

Elementary materials by the best of

$$r = A_c$$
 and (4, 0)

every other  $r_1^* = 0$ 

We get

DA by martiplying chiefs row IA state in partial management

AD .. . . . column A .. . . . . . . . . .

B, Ash by rw lit are programmed be rever y control

AE , V Type ammedit has long be sumb ( ) by (), (A) )

We consider especies vish as thing he matrices

$$\mathbf{U} = ((a\sharp)) \tag{4.7}$$

for which the diagonal elements  $u_i^* \Rightarrow u_i$  are where the representation of the matrix  $U^{-1} = U^{-1}$ , we see the matrix  $U^{-1} = U^{-1}$  are  $u_i$ , are suppressed to be updated.  $U^{-1}$  too is a matrix of the type (4,7).

$$\mathrm{Ep}\{(\lambda) = \mathrm{E}_{x,y}(-\lambda)$$

is an elementary matrix. Hence I we so How much you A from the left as well as from the right and a matrices of the same type B is said to be congruent as A, and is dinner in

I had A B I A C it follows the A of the R A A A A A A Hence the later is congruent to A for a loss of the extent of this example ongoing to extend the extent to be a symptomic to the forther producers and needs A

[43] i.e., a became that the formula of S and  $h_1$ ,  $h_m$  be arbitrary elements of S. The elements

$$a_1b_1 + ... + a_nb_n$$
 (4, 9)

the country of the strain of t

We have the most restrict submorate generate less the determinants

$$\Delta_{A+2+\cdots+} \Delta_{A+m_A} \tag{4, 10}$$

had a figure hafte mark A and let the chancets of A be not fit to a 2 th mark 1, the mark 1, the mark 1 and for every had be the ballest common to be for a common for the first common for the first common for a common for

$$\Lambda_{k-1} = \sigma_1^* \Lambda_{k-1+k_1} + \dots + \sigma_k^* \Lambda_{k-1-k_k}$$

horne to Marc. Thence in the sequence of modulo

Hence to Decrease I is a common be represented by the

$$\delta_{3} = e_{1}$$

$$\delta_{3} = e_{1}$$

$$\delta_{4} = e_{1}$$
(4. 11)



Theo am Contract restriction have the same alamentary divisors,

Proof the case of the term the left of referent the right solet with a pate x of an analytic the right x of x of

$$\mathbb{E}_{x,t}(\lambda | \mathbf{A} \longrightarrow \mathbb{E}_{x,t}(-\lambda) \mid \mathbb{E}_{x,t}(\lambda) \mathbf{A}) = \mathbf{A}.$$

Ly the 1-p of zn - 1 on which treat to be und U the [4/4] following operation in a six of the later.

d full numeral fact that the state of the st

(a) (a) 
$$= \{(a) + (\beta_i) = (-\beta) \}$$
 (-16)

(f) 
$$(a) + (f) = (a) + (f)$$
  $(a)$ 

(I Sweep out fire race out by a concentration of the errors way the high different

$$a, ab_0, \dots, ab_n \longrightarrow a, 0, \dots = 0$$

we can after their with column 1 the non-section of the section to will appear with the property

$$N(b) < N(a_s)$$
 for every  $t$ .

Proof full at be an expent for which the trial that there are elements a and b for which

$$a_0 + ca_2 = b_1 \quad N(b) < N(a_1) \le N(a_2)$$
 (4, 12)

holds.



## ALGEBRA

By one min will be see can regisce an by his

If he a the he is the periods of the authority of a complete field and the matrix, we can receive hy come on a compact that the matrix will appear in the agency for which are a compact to an appearance of the agency for which are a compact to an appearance of the agency for which are a compact to a comp

from 1 z is and a 'f we that we can express to be the first exhibits to be and the first exhibits to be a war to be a trade or reation the A of the elements of the matrix will be be a trade to be a trade or reation as earn of the not divisible by a fact and be as affected by the coverand orders to a fact and be as affected by the coverand orders to a fact and be appeared by the appearance of the matrix with the postern as many the appearance.

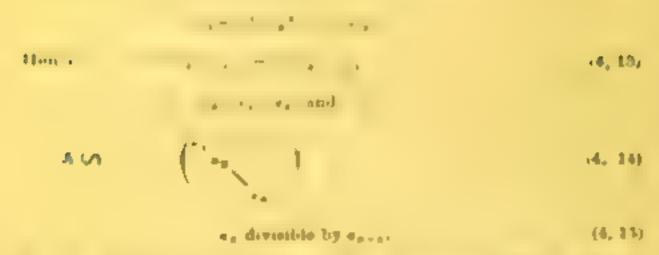
With the hip of this a pecation the "exception on analytic dime that a his un even at if the matrix of a which that function V has a man main and if a confidence from the hint a, of the elements of A then we can after A in any his anisometric that there appears an element high the function of taken integral positive walker only hinder after a finite number of steps the propositive must step, and that is only possible fan accurate of the matrix becomes equal to a.

By 1) e, went a proced on the left apper permer of the quatrix, and by the help of 2, the tree row and the dreep again white except out. So we get

Exhibit months to a fact the many that the form a country to a few periods are the get therefore



The h. e. f. of the minors of degree h is



If it the office out a to 141 to the feter court france is the first the first time of the first time for th

The or Perce makes A with a messa from S a regrete to characteristic to 14 as a first elementary are a strong to the absence of A. I force a there for an 1, is arrapped to a contrast the same force of a contrast to the same force of a same in a same the same in a same the same in a same to the same in a same in a same to the same in a sam

respective to a long on orders to a this was appropriate to a sure of the same appropriate to the same appropriate to the same and the same appropriate to the same and the same and the same appropriate to the same and the same and the same and the same and the same are as the same and the s

Theorem. A Ut B if and only if

where s' C, were matrices in a classicate from S, their determinants being unities of S

In transfermation of a victor space by the matrix Alleg ven by

Int 
$$C_1(x) = (y)_x \cdot (x) = C_1^{-1}(y),$$
 
$$C_2 \cdot \omega \cdot (x) = C_2^{-1}(y),$$

$$(y) = D(a)$$
 (4, 17)

Lar ets

Fig. 8. Let  $u_{-}u_{+}$ ,  $v_{-}$ ) in top n text vectors in the Euclidean appear, and  $v_{-}$ ,  $v_{-}$ , and great numbers. If the we turn  $v_{-}$  to  $v_{-}$ ,  $v_{-}$ 

[4/6] We shall now apply the theory for open many explicit from a ring S not to an open in robot on the factor of the theory for the constitution of the theory for the constitution of the terms from the northward form got now will then appear in a new light.

the contract of the later to the property of the property of the later to the property of the

In the tends of the part is a factor of the normal form of the normal formal formal



following congruences

$$\begin{pmatrix} a & 1 \\ 0 & a \end{pmatrix} \wedge \begin{pmatrix} a^{n-1} - a^{n} \\ 1 & a \end{pmatrix} \wedge \begin{pmatrix} 1 & a^{n-1} \\ 1 & a^{n-1} \end{pmatrix} \wedge \begin{pmatrix} 1 & a^{n$$

and, for (a, b) = au + bu = 1,

$$\left(\begin{array}{c} z & 0 \\ 1 & 0 \end{array}\right) \wedge \left(\begin{array}{c} z & 0 \\ 1 & 1 \end{array}\right) \wedge \left(\begin{array}{c} 1 \\ 1 & 1 \end{array}\right) \wedge \left(\begin{array}{c} 0 & 1 \\ 1 & 0 \end{array}\right) \wedge \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}\right) \rightarrow \left(\begin{array}{c} 3 & 0 \\ 1 & 0 \end{array}$$

We take the transfer of the space of the plant of the space of the spa

· button and share near t 
$$\begin{pmatrix} 1 & 0 \\ 0 & (\lambda_x - e) \end{pmatrix}$$
 the attent is not being

white of the control the control of the common the control of the common the control of the common that the common the control of the common that the common the control of the common that th

1. 
$$1 - c_1 - x_1^{\alpha_1}$$
  
1.  $1, (A_1 - x_1^{\alpha_1})$   
(4. 21)

$$1 = 1 + 1 + 1 + 1 = 1$$
where  $+$ ,  $\geq 2$ ,  $+$   $= 2$ ,  $+$   $+$   $= 1$ . (1 =2)

or la, bus the sum of some form as in the mount of the decimal of the sum of the first of the sum of the first of the sum of the first of the sum of the s

1, ..., 1, 
$$(\lambda_1 + z)^{\theta_1}$$
, ...,  $(\lambda_1 + z)^{\theta_2}$ ,  $(\lambda_1 + z)^{\theta_2^{\theta_2}}$ . (4, 23)

From all 32 of the politic 4, 23 is the gorman from 1 (1) the of the later tensor that the the resemble sentence of A we can reduce the observer the observer the observer that the observer the product to a normal broad \$1.2. So we also be found in a recommendation to A the designation with being powers of a common and applying the congruence in the 120 or get a topomorphic

$$1, \dots, 1, \psi^{\pm}(x), \dots, \psi^{\pm}(x), \psi(x)$$
, (4, 24)

where y relief to the auto press of a the free detects the precise of the precise of the free detects. Hence the first the precise from 4.14 of A. 21. If (4.24) is a settle we can find out the rests A. and the expension e. if a factor of the free detects the free from A. 2.) a compact of the classes of quatrices A. free detects the weight the compact of a test the expension to transfer into the electric of a test the consists of the root A. (2) and congruent to A. 21. There is an if I arrespondence that the electric the electric the contract of the root A. (2) and congruent to A. 21. There is an if I arrespondence that the contract of the contrac

[6/1] Let K be a sub-field of 5 and [5: X]=2

The property of the property o

\* an account on it a the enter the people to of two elements a terral number of the extrapolating enter the following the follow

See to

$$\hat{A} = A \qquad (0, 2)$$

.



the the preparation of the same of the context 1,2; approximation of the same of the context 1,2;

For both cases we make the forming and in-

from (0.0) them there exists an overer y 2 , a hitter

$$\omega + \beta \beta = \gamma \gamma. \tag{5, 0}$$

This condition is satisfied a g

1 If he has the he dof the real namers

2 of K to the first fither each market to the fither maples made

the thorny has mady been a 1 th there are to a factor and percent that are to the term and to be the first and the form a gette property of the manifestation specifically the first and the form a first and the form and the first and the fir

$$a_1 a_2 + \dots + a_n a_n = AA$$
. (5, 30)

As to become to K from s = ) it is no that a first it is fall factors, and therefore both at u i be equal to 0 for each it was if the organization (5 % and 1 ) are threefore believed from lifer each (a fighter) as any part former to 5 km f a which each part of elements of K.

$$\frac{a_0q^{k_0}-1}{a_0q^{k_0}-1}, \quad a^{q_0}q^{k_0}-1. \tag{6.4}$$
 expectably 
$$a_0q^{k_0}-1, \quad a^{q_0}q^{k_0}-1.$$

As the promote that the same of the transfer o

Therein Lettel stranger can a state of a

102

ALGERSA

conditions.

$$\frac{2t}{t} + \left( \frac{1}{t} + \frac{1}{t} \right) = 1 \tag{5.6}$$

and ....

$$A_1u! = e!$$
 for  $b = 1, ..., n$ . (5, 7)

P + 1 , which tests notices  $0 = 1 + 6 + \frac{5}{6} + \frac{1}{12}$ ,  $0 = \frac{5}{6} + \frac{1}{12} +$ 

s den , f  $\geq r^2$  ,  $\geq r^2$  = then  $\geq r$ ,  $r_s^2$  = bolds, for  $r^2$  2 +  $r^2$  ,  $r_s$  . We can solution for the solution the set get  $r^2$  as the  $r_s$  as aform the set by

 $\lambda = \lambda_1 \lambda_2 + \lambda_3 \beta$  (see the  $\alpha = 11$  no.,  $\beta_1 \lambda_2$ ), alterior the conditions (5, 5), (5, 6) and (5, 7).

From 0 for which K = 2 and  $C_1 \ge 1$  hada

[5/8] Proceedit astiff a divide, 2) can be considered as a proporty of the most to, the superior test of the kind in a simple form at a begin to use the fillering notations.

Let A ... we have a with elements from A then

From these formulas it follows that

$$(AB) = (A)^{\dagger}$$
  $AB = AB$   $AI)^{\dagger} = B^{\dagger} A^{\dagger}$ 

and for det Ask0,

$$A_{2}^{-1} = (1)^{2} A^{-1} = 1^{2} (1)^{2} (1)^{2} = A^{-1}^{2}$$
 (5.11)



The equations , and for an observiors open out with

$$((u_1^*))^{\frac{1}{4}} \approx ((u_1^*))^{-4},$$
 (5, 12)

A zintera an information of the control and the beautiful An the transpost forms (or a ) matrix is also instant a six and are six and also antiarly

$$\Sigma u_{1}^{*}u_{1}^{*}w_{0}$$
, for  $i\neq j$  (6, 64)

$$\sum_{i} u_i^* u_i^* = 1, \qquad (5, 47)$$

As dot  $((u_i^*))$  mdet  $((u_i^*))^*$  mdet  $((u_i^*)^{-1} - 1 + \det ((u_i^*))$  holds, hence

$$\det ((0!)) = 1,$$
 (5, 18)

A andary matrix remains sustain after an ethnic permuta in I the rows and after summer the poduct of unitary material anctary

$$U = \begin{pmatrix} \pm 1 \\ \hline U' \end{pmatrix} \qquad (5.14)$$

and most the matrices to the unitary that here a sixter) to the matrix of a sixter and the termination of the large the sixter and the sixter and the termination of the theorem to the the theorem to th

Theorem locary vector is specially as an arrangement of the property that a given two points. Here is a fact of only.

If the content of a matrix II at the first to the polynomial belong to A and

$$\mathbf{H} = \mathbf{H}^{\dagger} \tag{6.15}$$

H to said to be an Harmitian matrix.

but II as an If constant matrix U a on large, frix (b) o II; U 11), has the same characteristic polynomia, as H and H<sup>2</sup><sub>1</sub>, C<sup>†</sup> H<sup>2</sup><sub>1</sub>, Mapper H<sub>2</sub> is Hermitian. Let 3 be specifically contained by an in Fig. (a).

<sup>\*</sup> This could from in examples of an entradic balance. It is a closed bard, e.g. t. e. first of the complex numbers

a vocal for which (H-V),  $\rho t=0$  halfs. Such, a vert  $(\beta)$  must just an 2-1 and 1+1 the number 1+1 must expend that 1+1 the last number 1+1 transferred in the same means 1+1 that is the first number 1+1 the another 1+1 the another 1+1 the another 1+1 the another 1+1 the same means 1+1 that 1+1 the first that the shift 1+1 that 1+1 the first 1+1 the another 1+1 that 1+1 the same 1+1 that 1+

$$H_1 = \begin{pmatrix} H \end{pmatrix}$$

Hence the an executed K. H. with Hermitian matrix of degree n = -W can transfer to by a minury matrix of the same manner as H. has been transfermed.

$$U + U = \begin{pmatrix} A' \\ -11' \end{pmatrix}$$
 As the instrict  $U_3 = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$  by

niso winterly 
$$U_1 H_1 U_2 = \begin{pmatrix} \lambda & & \\ & & \end{pmatrix}$$
 or an Hermitian matrix  $H^{\mu\nu}$ 

After his the weight H. constrains limited log no matrix by a matrix abids of the product. I no tary matrix hand a therefore unitary. So we get the following superstant theorem.

The rem. An Hermitan matrix H can be transformed by a unitary metric attraction to a dog on main and the risks of gatz belong to be

Let k be the first fibercal now bere the the field of the complex numbers, then it files from the their in that the reta of the character is a particle of the character is a particle of the character is a symmetric real mater. He was an ensuing me mater with real clements for which a, it has not from the a, where K the against the field of the remarkers to the material the character stip polynomial being to A Hence we can apply the present, the remark on this case, the undury matrices becoming now orthogonal matrices and we get the

Can very If in a matrix A with real numbers as extremts at =a; holds, then  $\chi_1(x)$ , has real mosts may and A can be transformed by an arthogonal matrix with real coefficients into a diagonal matrix.

The theory of motives will now be applied to a near and quadratic [6/4]



forms. We introduce a flat cost of and fin see

$$\Sigma = \Lambda \left[ x_1, \dots, x_1, \dots \right] \tag{6, 17}$$

of the personness to the indepetes (5.10 with coefficients from A. This extension will be made so that every indefine (5.16) becomes replaced by the conjugate. So we have samply to report an very payment leach north contains and each independent by the employate. In the case K = A the automorphism, in the dentity. In every case there belongs to every element configurations of A. magnety defined conjugate changes countries.

In 1,3, a vector has been represented see (1-25) by a matrix, the heat column of which is formed by the coordinates the alements of the other colon as being 0. We will now apply the notations of conjugate matrix and transposed matrix to these special matrices, so that

$$(z) = \begin{pmatrix} z_1 & 0 & 0 \\ \vdots & \vdots & \ddots & 0 \\ z_n & 0 & 0 \end{pmatrix}, \quad z) = \begin{pmatrix} z_1 & 0 & 0 \\ \vdots & \vdots & \ddots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \vdots & \ddots & \vdots \end{pmatrix}$$

$$5, 18)$$

$$5, 18)$$

Let  $\Lambda = \{f_{i}\}$  i, then A(g) is a matrix with the first element

$$\sum a_{1}^{2}x_{1}y_{0}$$
 (5, 19)

all other elements being equal to 0. To every metrix A. there corresponds

a Character 19 and exercely bet

$$(x)=B(x'), y=C(y'), not$$

$$B^{\bullet} AC = A' = ((a')),$$
 (5, 20)

then

$$(x')^{\circ} B^{\circ} AC(y') = (x)^{\circ} A(y).$$

Hence 
$$\sum v_1^* x_1 y_2 = \sum a_1^* x_1^* x_2 x_3$$
 (5, 20)

The fermulae (5, 20, and c) 20° g we the transformation of b linear forms

We will now consider the case where A is an Herin tion matrix, and where s and y are conjugate. Then It and C are as a conjugate, or, the becomes an Hermitian form

$$\Sigma a_1^{\dagger} x_1 \overline{x}_4$$
, where  $a_1^{\dagger} = a_1^{\dagger}$ , (5, 21)

and 
$$\sum_{i} x_i x_i = \sum_{i} (x_i x_i)$$
, where  $\sigma(1) = A = C + AC$  (5.22)

A is the refere an Herm t an instruction. In the age of non-work K=1 at n = n and  $x_1 = x_1$ , the bilinear from becomes quadrate and the Herm an matrix becomes a symmetric on:

$$\sum_{i=1}^{n} x_i x_i = \sum_{i=1}^{n} x_i x_i + 2\sum_{i=1}^{n} x_i x_i$$
 (5, 28)

If the characteristic of K is different from 2 every quadratic form in  $x_1 - x_2$  can be represented in the form 5-25. We will omit the case of histories 2 which needs a special treatment limited there is an 1,3, correspondence between the quadratic ferms and the symmetry quadratics. The transformation is done by

$$x_{i,j}(x) = X_i = C^{(4)}A^{(4)}$$
  $x_i = C^{(4)}x_i + \sum_i x_i^{\dagger}x_i + \sum_i a_i (x_i, x_i)$  (5. 24)

The formula for the transformation of quides to firm remark the of the transformation functions [see 1 to ] at the notion of inverse matrix has been replaced here by the note in formula and that a trace any that did to On applying the through of [5 5] and its encounty to (4 22) and 5, 24) we get the following for lamental theorem.

The rem Every Hermit so form one be transformed by a unitary transformation into the normal form

$$\Sigma_{0,\pi,\Xi_{1}}$$
. (5.26)



and every quadrat. form with coefficients from hican be transformed by an orthogonal transfer in it on with coefficients from K ato the normal form.

$$\sum_{i=1}^{n} x_{i}^{-1}$$
 (5, 26)

a, being conents of K not the cases

Let K = X be the first of the real numbers, then every quadrate form [6,5] can be transferred as

$$Q(x) = a_1 x_1^{-1} + \dots + a_r x_r^{-1} \qquad (5, 27)$$

where the efficiency of the restriction area by the transformation with an other material to the determinant for other noticed by the transformation and the first and the restriction and a therefore a post to a life we repose a try of the formula [5, 27) and not a start for a control to the sign of the attention transformation because him, then we can transform an artiferation with determinant of the restriction and the same transformation with determinant to the same transformation in the direction of the same transformation in well him which the region. Then transformation is well him which the region that transformation is well him which the region the principal exest.

Let  $a_1 = a b_1^{-2}$ , and  $b_1 x_1 = a_1$  then  $(b_1 x_1)$  is transformed to a sum of equation with cortain  $a_{00} a_1 + a_2$ . After a permutation of the induction of the i

$$q(x) = x_1^2 + ... + x_n^2 - x_{n-1}^2 - ... - x_n^2, (6, 28)$$

Hence every quadratic forms on be transformed to all 29 by a non-legenerated linear transferre at a with real coefficients. The integer r a therefore of the quadratic form and therefore measured: We was prove that p is an invariant too.

The rem living quadrates from with real conficients can be trained into 1 into the action y such must form, ribbar non-deginers ed linear transformation with real sould length

From a always possible and reson average. We have therefore a y to prove that a transformation of q(x) into

$$q = -\frac{1}{2} + -\frac{1}{2} + \frac{1}{2} - -\frac{1}{2}$$

by the ar non-degenerated substitution is possible only if p=q. Let  $p \nmid q$ , any  $p \geq q$  without any loss of generality.

$$x_i = b\{x_1 + ... + b\}x_i$$
  
 $a_i = a\{x_1 + ... + a\}x_i$   $i=1, ..., \tau$ 

 $q(x) = q_1(x)$  for corresponding existence  $(x_1, \dots, x_r)$  and  $(x_1, \dots, x_r)$ . The q+r-p < r mean homogeneous equations

$$c_1^*x_1 + \cdots + c_r^*x_r = 0$$
  $k=1, ..., q$ 
 $x_r = 0$   $t = p+1$   $t$ 

become solution  $\xi_1$ ,  $\xi_2$ ,  $\theta_1$ ,  $\theta_2$ ,  $\theta_3$  different from  $(0, \dots, 0)$  one the rank of the matrix of the system if equations is  $\leq q+r-p < r$ . The e-groupe of this values of  $x_1$ ,  $x_2$ , are  $\theta_1$ ,  $\theta_2$ ,  $\theta_3$ ,  $\theta_4$ ,

A quadratic form is and to be positive definite if n=r=p it is negative definite if n=r, p=0, it is seem definite if n>r, p=r=r=0, and it is indefinite if r>p>0

[5/8] Finally we will give a geometrical interpretation of the cast results without going into the details

In the projective in - 1, d.mensional space, a quadric becomes repre-

$$\rho \Sigma a_1^* x_1 x_2 = 0$$
,  
where  $\rho \stackrel{*}{=} 0$  is an arbitrary factor.

Honce the quair c has one and only one acrina form

$$\rho q_1(x) = 0,$$

and the sign of a can be fixed in such a manner that q(r) has not fewer positive than negative terms. We get therefore the different types of quadres in the projective in -1 dimensional space given by the different normal-forms

$$q(z) = 0$$
, for  $t=1$ ,  $n / + \le p \le t$  (5, 20)

Especially the quadries without any real point are these for which permit. The nominal form 5.20 has the property that every fundamental point of the coord-nate system or every point, for which all the coord-nates are equal to 0 except  $r_t$ , whilst to the opposite hyperpoints  $r_t = 0$ .



In the affine a d trengional space the quadrics are given by

$$\Sigma a_1^* x_1 x_2 + \Sigma b_1 x_1 + a = 0.$$

On applying the theorem of [ 5] we get easily the following types affine normal forms.

$$q(x) = 0,$$
  $r = 1, ..., n$   $r/2 \le p \le r$   $q(x) = 1,$   $r = 1, ..., n,$   $p \le r$  (5, 20)  $q(x) = a_{r+1},$   $r = 1, ..., n-1,$   $r/2 \le p \le n$ 

If we replace que by the quadratic form Q zerof 5, 27) we get the types of quadrics different in the sense of matric geometry. The formula 5, 27) can use be interpreted for the jet by generally.

Let R z<sub>1</sub> = z<sub>n</sub>) and S /z = z<sub>n</sub> = t<sub>n</sub> · q indicate forms. S representing a quadric without rest points thin we transform be conducted in a noh is manner that S is transformed to S = z<sub>n</sub><sup>2</sup> + + +z<sub>n</sub><sup>2</sup>, and R to R'. By any arthogonal transformation S without he altered but we can transform R + y an arth good transformation to the normal-form the Z?

Hence we can transform significant species y

The e ements of the the punch have therefore term transfermed to the normal-form simultaneously

## 5 6. RESULTANTO

Let

[6/1]

$$f(x) = x^n + a_1 x^{n-1} + \cdots + a_n = (x - x_1) - (x - x_n)$$
 (6.4)

$$g(x) = x^{n} + b_{n}x^{n-1} + \cdots + b_{n} = (x - \beta_{n}) \quad (x - \beta_{n}) \quad (0, 2)$$

$$R(f,g) = \prod_{i} \prod_{\alpha} \{\alpha_{i} - \beta_{i}\}, \qquad (6,8)$$

Then R(f, j) is un prety defined by the plane do f and j and j and if a said to be the exact set of f and j. The necessary and authorist condition that that f and j is ay have a common root a that the resultant is equal to zero

From (6, 8) it follows

$$R(f, g) = (-1)^{mn}R(g, f);$$
 (6, 4)

from (6, 2) and (6, 8)

$$R(f, g) = \prod_{i=1}^{n} g(a_i),$$
 (6, 6)

and by into rehanging family and approving 6 4 we get

$$R(f, g) = (-1)^{mn} \prod_{k=1}^{m} f(\beta_k), \qquad (6, 6)$$

The right sit of the cips in the a symmetric parameter p anomal in  $a_1$ ,  $a_2$  with the  $a_1$   $b_2$   $b_3$ ,  $b_4$  there p lyn main h having atomic coefficients. From Part II [1, b] it follows the R p pumber represents 1 as a p-vision and in the occurrentary symmetric p-lynomials of  $a_1$ ,  $a_2$  with such case  $p_1 = p_2$   $h_1$ , h-visions by a physical of  $a_1$  and  $a_2$  are the coefficient h the coefficient h in the coefficient h is a polynomial of h and h are h and h are h and h are h are h and h are h are h and h are h and h are h are h and h are h are h and h are h and h are h are h are h and h are h are h are h and h are h are h are h are h are h are h and h are h and h are h are h are h are h are h and h are h are h are h are h are h and h are h are h are h and h are h are h and h are h are h are h are h and h are h are h are h are h and h are h are h are h and h are h are h are h are h are h and h are h are h are h are h and h are h are h are h and h are h are h are h and h are h are h are h are h and h are h are h and h are h are h and h are h are h are h are h and h are h are h and h are h are h are h are h and h are h are h are h are h are h and h are h and h are h and h are h and h are h

by a common which will be sent to be to the representation will be written as

$$B_i(f, g) = R_i(1, a_1, ..., a_n; 1, b_1, ..., b_n)$$
, (6, 7)

If in any term  $A = a_1 - a_2 = a_n - b_1 - b_2 = b_n$  the factors  $a_0$  are

represent 145 . . . , and the fact ra 4, by  $\beta_1 = \beta_n$  then A becomes an homogenious polynomial in  $\alpha$ , and  $\beta$ , of degree \*

$$t = x_1 + 2x_2 + \dots + x_n + t_1 + 2t_2 + \dots + mx_n$$
 (6, 6)

is easily to the couple of a form h is it follows that H(t,y) is in general a flagrance between the term of H is  $a_1$ ,  $a_2$ ,  $b_3$ ,  $b_4$ ,

Second of the contraction of the second the property that Second of the contract of the property that the second of the second o



th Bi, Ma we get

\*\*\*\*

$$S = \Sigma = \Sigma (\alpha_1, ..., \alpha_n, \beta_1, ..., \beta_n)$$

The right side of this equation is equal to  $x \in \mathcal{A}$  of  $\alpha_1 = \beta_1$ Hence

$$\Sigma = \Sigma (a_1, \dots, a_n, \beta_1), \quad \beta_n = \Sigma (a_1 \dots a_n, a_1, \dots, \beta_n)$$

Bubtracting the corresponding terms in heir ght ade we see that 2 is distable by ap \$\beta\_1\$ and a the rear member ? I have loss \$\beta\_1\$ to dward a by a, B, hence X a divaster by R tr Or I (1 ) 2 It ( ) From the 2 theorem of Pet II , 10 , at f we observe that  $R(1, a_1, ..., a_n; 1, b_1, ..., b_n) = (B, R(1, a_1, ..., a_n, 1, b_1, ..., b_n)$ a I vooble by the result out. The weight of every term is therefore in t less than in n, if each term has the weight a n is defice from the resultant by a factor of weight 0 only and the factor is the coefficient of the term b. am S.

To get the resultant as a polynomial in ", a, b, a, base [6/2] therefore to find out a polynomie B, with the following three properties

- b ≈ 0 if f and g have a common root
- 2 Luch term of S has the weight o n
- 3. The term b\_" has the coefficient 1

A po yourmal of the kind can car ly be found out by the following consideration

Let f w, =0 = g a then the following n+m equations hall

$$a^{m-1}/a = a^{m-1} + a_{m-1}a^m + a_na^{m-1}$$
 (6)

$$n = f_{-1} = \frac{1}{n^{n+1} + 1} \frac{1}{n^n + \dots + n} = 0$$

$$a^{-1} g(a) = a^{-1} + b_1 a^{-1} + b_{-1}^{-1} + b_{-1}^{-1}$$

$$a^{\frac{1}{12}} \cdot y(a) = a^{-1} + a^{-1} + b_{-1} \cdot a^{-1} + b_{-2} \cdot a^{-1}$$
 (20)

$$a^{n+1} + b_1 u^n + \dots + b_{n^n} \qquad b$$

We consider this system as a system of neutropy strong in a train, and it can be satisfied unly first differentiate a equal to zero. Hence a necessary condition for that fixed y may have a common root or

The terms of S are of weight one can't the term but is the diagonal element and has the coefficient 1. Hence

$$S = R (f, g), \qquad (6, 9)$$
Let  $F(x) = a_x x^n + ... + a_n$ 

$$O(x) = b_x x^n + ... + b_n.$$

[6,8]

white pothing is surpered about the coefficients. We define now

$$R(F, G) = R(a_1, ..., a_n, b_1, ..., b_n)$$

$$a_1 ..., a_n$$

$$a_n ..., a_n ..., a_n ..., (6, 10)$$

As in 6 if there are a rows with elements a, and a rows with elements  $b_1$ . If  $a_2 = b_2 = 1$ , then F = f, G = g, and we see from (0, 0) and (0, 0), that the notation f(F, G) = 0 is not to f(F, G) = 0. We have to consider three doses.

1. 
$$a_n \neq 0$$
,  $b_n \neq 0$ ,
$$F(x) = a_x \left( x^n + \cdots + \frac{a_n}{a_n} \right) = a_x \phi(x_1 + a_1) \cdots (x_n + a_n)$$

$$C_1(x) = b_x \left( x^n + \cdots + \frac{b_n}{b_n} \right) = b_x \psi(x_1 + b_1) \cdots (x_n + b_n)$$



From (8,10) it follows that"

$$R(F, G) = a^{\alpha}b^{\alpha}_{\beta}R(\phi, \phi).$$
 (6, 11)

Hence for (6, 3) (6, 5), and (6, 6)

$$R(F_{x}|G) = a^{\alpha \beta} \prod_{i=1}^{n} \prod_{\alpha_{i} = \alpha_{i}} \alpha_{i} = \alpha_{i} = -2\Pi(I_{\alpha_{i}}) = (-1)^{\alpha_{i} + \beta_{i} + \prod_{i} I_{i} + \beta_{i}'})$$
  $0', 12$ 

Hence in this case R/F G = 0 a the necessary and sutherent condit in fire

 $(2-\alpha_{\mu}=0)$  by (0) From 6 to 4 f we the H  $(1,\alpha)=0$ , reference of a small street

$$b_s = b_s + 0$$
,  $b_s = 0$  (or  $a_s = 0$ ,  $b_s + 0$ ).

Let  $h_1 = -mh_n = 0$ , then R is (n) = 0 and every root of F entirelies obviously G(n') = 0.

Let every 
$$h_{4 < \tau} = 0$$
  $h_{-} \neq 0$   $h_{+} \neq 0$   $h_{+} x = h_{+} x^{n-1} + \cdots + h_{n}$  for  $=1$ 

By ketting 0 for h. in t. 10 we get R.F. tr = a. RtF, G. v.

House Ref. if will dash early than 141, have eccommon took or if F and O have a sum out rate his excreasionding to particular had for a particular to the first on the first on the first on the first one of the

Theorem If F and G i have a proportion than R F (1 = 0). If 1. F, G = 0. then either a = b = 0. The and the base is computed.

Exercises 1 Consider the service of a 12 and 0 12 a the case 3.

State the necessary and out that on him for hite # "

Let  $u_1$ ,  $u_{m+1} = u_1$  be the classical of the first example of the  $\{u, u\}$  determinant (0, 10), and let

When we multiply the n + m equations

$$g^{n-1}F(s) = a_n x^{n+m-1} + \dots + a_n x^{n+1}$$
 
$$F(s) = a_n x^n + \dots + a_n$$
 
$$x^{n-1}G(s) = b_n x^{n+m-1} + \dots + b_m x^{n-1}$$
 
$$G(s) = b_n x^m + \dots + b_m$$

with ut. .... une with me respectively and add, we get

$$w(x) F(x) + v(x) G(x) = R(F, G).$$
 (6, 15)

We can express this formula by the following theorem.

Theorem. Let R be the ring generated by the indefinite x and the coefficients of F and G, then R(F, G) is linearly dependent on F and G with coefficients from R.

Exercise. Prove the theorem of [6/3], without any reference to symmetric functions, by the help of the last theorem. (Special attention should be given to the case when every noisotor is equal to zero.)

If F(s) and G(x) have a common root, the highest common factor (F, G) is a polynomial of positive degree; we can get it by the algorithmus of the L.c.f., hence its coefficients belong to the ring generated by the coefficients of F and G by addition, subtraction and multiplication.

[6/5] Let the coefficients a<sub>x</sub>, ..., a<sub>x</sub>; b<sub>x</sub>, ..., b<sub>x</sub> of F(x) and G(x) be polynomials of K(x). K being an arbitrary field;

$$F(s) = f(s, y), \ O(s) = g(s, y).$$

The resultant R(F, G) is therefore a polynomial in y, and from (6, 13)

$$R(F, G) = R(y) = u(x, y) f(x, y) * v(x, y) g(x, y).$$
 (6, 14)

We suppose that at least one of the polynomials  $a_* = a_*(y)$ ,  $b_* = b_*(y)$  is different from the polynomial 0. From  $\{6/3\}$  it follows that f(x,y) and g(x,y) have a common factor depending on x, if and only if R(y) is the polynomial 0. The procedure of getting R(y) and g(x,y) is called elimination of x. Let q be a root of R(y); then R(f(x,q),g(x,q),=R(q)=0. Hence either  $a_*(q)=b_*(q)=0$ , or f(x,q),g(x,q) have a common root. By this method we can find out the common solutions of two equations

$$f(x, y) = 0, \quad g(x, y) = 0,$$

## CORRECTIONS

Part I. (see the corrections given in Part II.)

Page.	Line.	Read	Por
ili (Proface)	10	does	do
1.0	19	(2/H)	(2)
16	24	depend on	apply to
20	10	10.	(11)
	28	11.	10.
Part II.			
.6	10	6+0	0+0
0	10	M	M*
11	22	(b'+ia)	(b'+t)
	23	n/so	uff
18	20	auto	and also the distri- butive law are
	22	a non-distributive system	a system
	28	the reader may verify that (2) generates an addition and multiplication of the classes which the commutative, a clatice, and distributive hald, and	i a prove that for mo-
18	21	field	fleld
18	21	0-1	0001
21	12 and 13 : interchange the exponents 1		
23	24	$f(x), \psi(x), \psi(x)$	$f(x), \psi(x), \psi(x)$
20	10	$\frac{a_n}{b_m}$	$\frac{b_{is}}{a_{n}}$
		84.	b <sub>n</sub>
	11	V <sub>ee</sub>	e <sub>a2</sub>
	26	$\phi(x) \phi_1(\chi)$	\$(x)\$1(x)
33	6	K	K <sub>1</sub>
	20	K(0)	K(a)
181	9.	N	K [three times  ]

## CORRECTIONS

Page.	Line,	Read	For.
42	1.0	to K(a)	to K
46	9	α <sub>θ</sub>	a <sub>f</sub>
47	19	21-1 + +b.	z==1+ +ba
	80	F(a1 , a4)	F(a1,, an)
52	17	Ä	b (error not in all copies)
61	8	(18, 2)	(13, 2
	5	(18, 2)	(13, 1)
	6	u(n-1)	n(n-1)

## Parts III ... V. (in this volume)

8	19	P	P [twice]
	.20	Q'	Q [twice]
	25	Reses	0:11-1
8	25	u - 8	a+1-a
11	li li	>6>	<e<< td=""></e<<>
19	4	+	- [between the fractions]
18	27	$\frac{P_{g,e}}{Q_{g,e}}$	Pun a
14	17	#x   A	Ani A
1.6	1	6	a
	15, 10	purely periodic	puzely
19	8	[8/4] [on the m	argin)
	. 9	>1	>0
20	12	α+β	α <sub>1</sub> +β <sub>1</sub>
		444	s+e; [bwice]
22	5	√26	28
81	1	Ser	24
	- 6	this sum is considered to the sum taken is divergent follows case Q <sub>2**1</sub> ->~. I	where it y that
	10	interchange "the	e" with "and"

Page.	Line.	aRead	For
41	17	(1, 3)	(B)
42	21	(=-a)(a-a)	(x-a)(x-a)
58	16	(1, 80)	(80)
58	6	9-40824	2*40224
ă0	15	$\left(1+\frac{2b_{3}^{-}+b_{3}^{-}}{b_{4}^{-}}\right)$	$\frac{(1+2b_3^m+b_4^m)}{b_1^m}$
61	14	b <sub>1</sub>	b <sub>4</sub>
66	15	[6/2]	[6/]
68	22	$\Delta_{2-0}$	$\Delta_{4i}$
72	7	(see Part II [1/2])	(see Part II [1/2].
	18	modul M	modul m
	30	$\{(e_i^*)\}$ , where $e_i^*=1$ , and $e_i^*$	(e), where a = 1, and a
74	7	exists	estint
75	18	al-	al:
76	21	exists	exist
70	26	[omit=]	
80	10	[place A from the 2nd to the line of (2, 14)]	e lat column in the lat
	11:	A <sup>(r)</sup>	A <sup>(n)</sup>
		X <sub>A</sub> (r)	XAISI
82	14	[replace the index m by q)	
85	G.	W3.1	Wiseel
101	17	existe	exist
105	28	(x)*A(y)	A(y)
100	100		3